

High-Dimensional Statistics & Sparsity

UDRC Summer School

João F. C. Mota

Heriot-Watt University, Edinburgh, UK



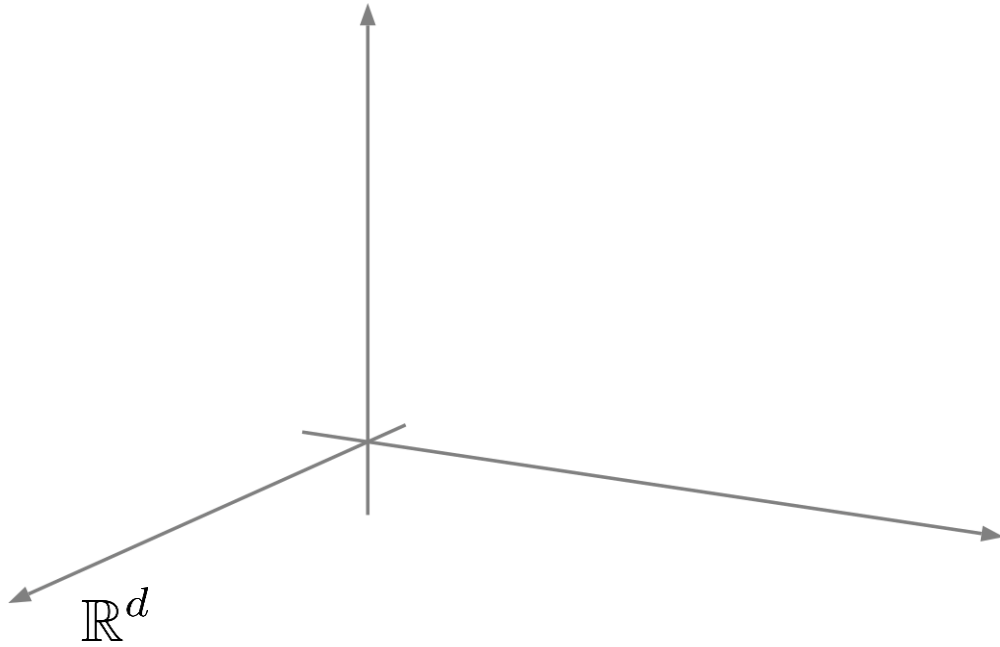
Motivation

Motivation

Hypothesis testing in high-dimensions

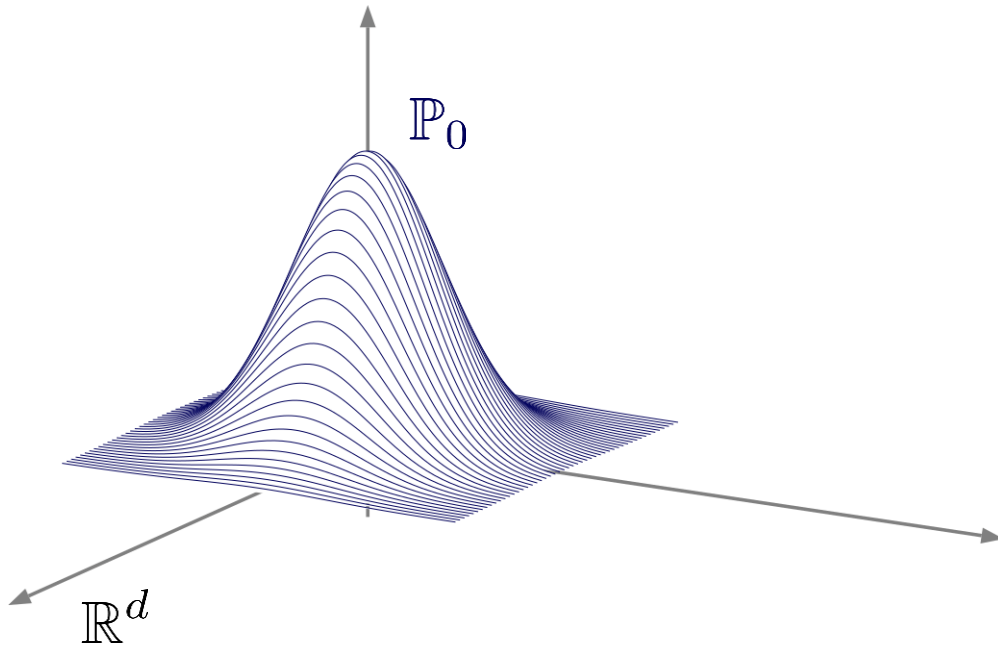
Motivation

Hypothesis testing in high-dimensions



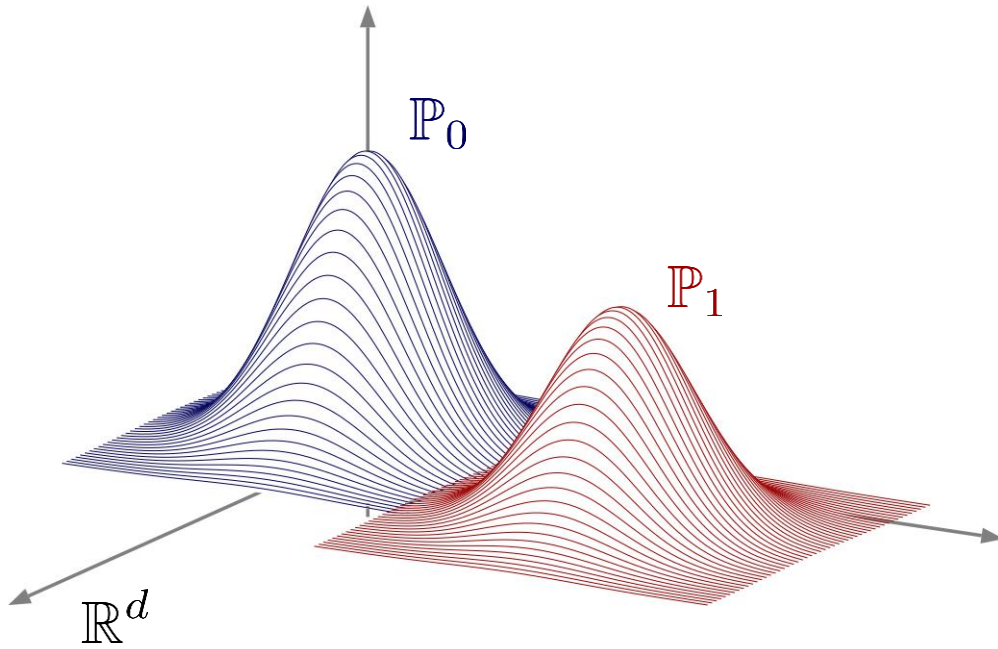
Motivation

Hypothesis testing in high-dimensions



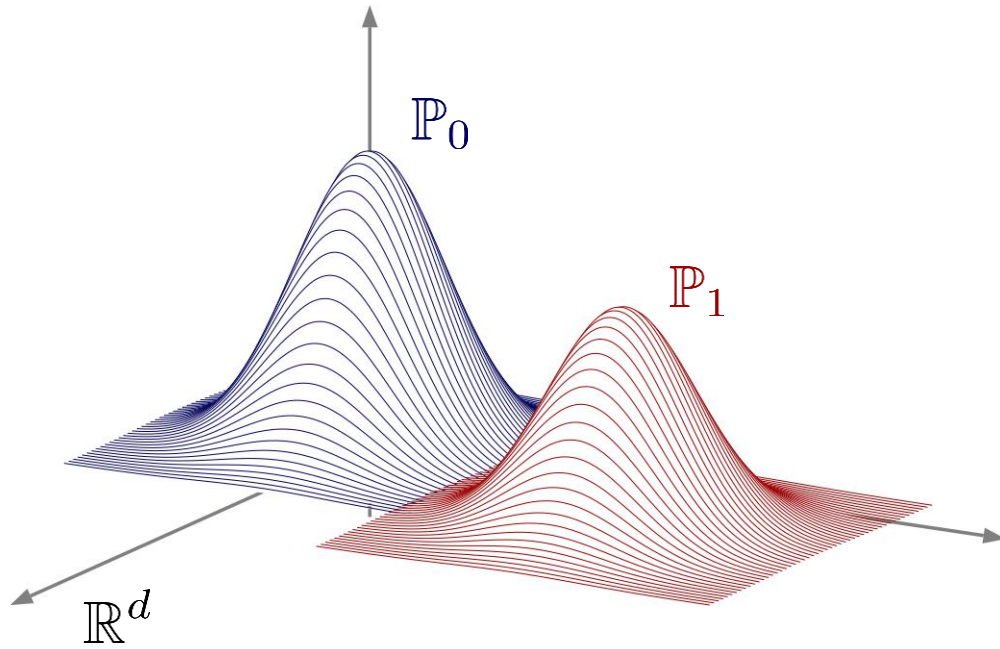
Motivation

Hypothesis testing in high-dimensions



Motivation

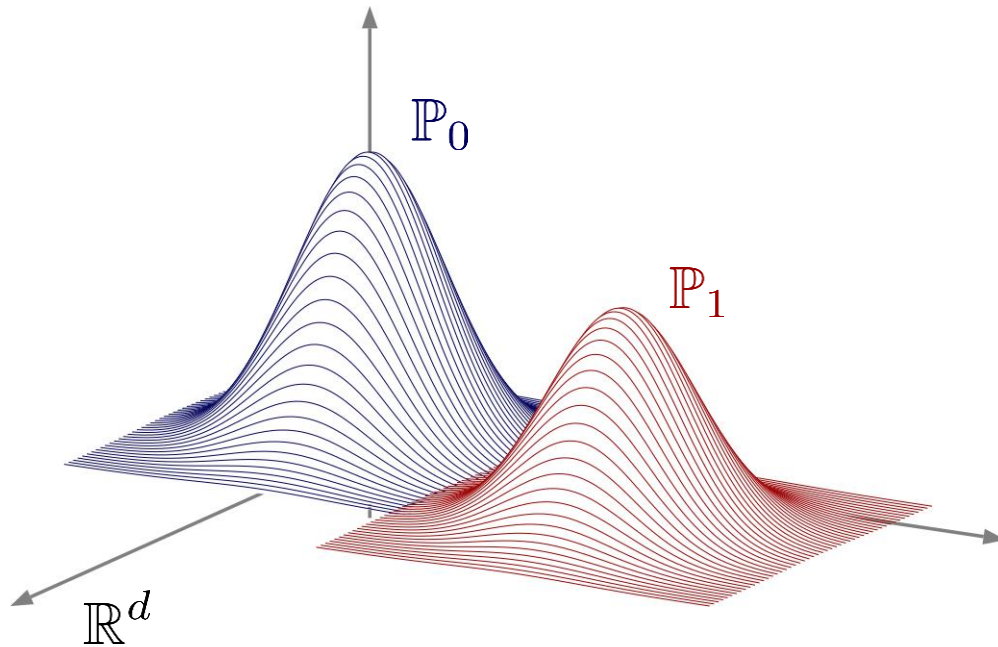
Hypothesis testing in high-dimensions



Problem

Motivation

Hypothesis testing in high-dimensions

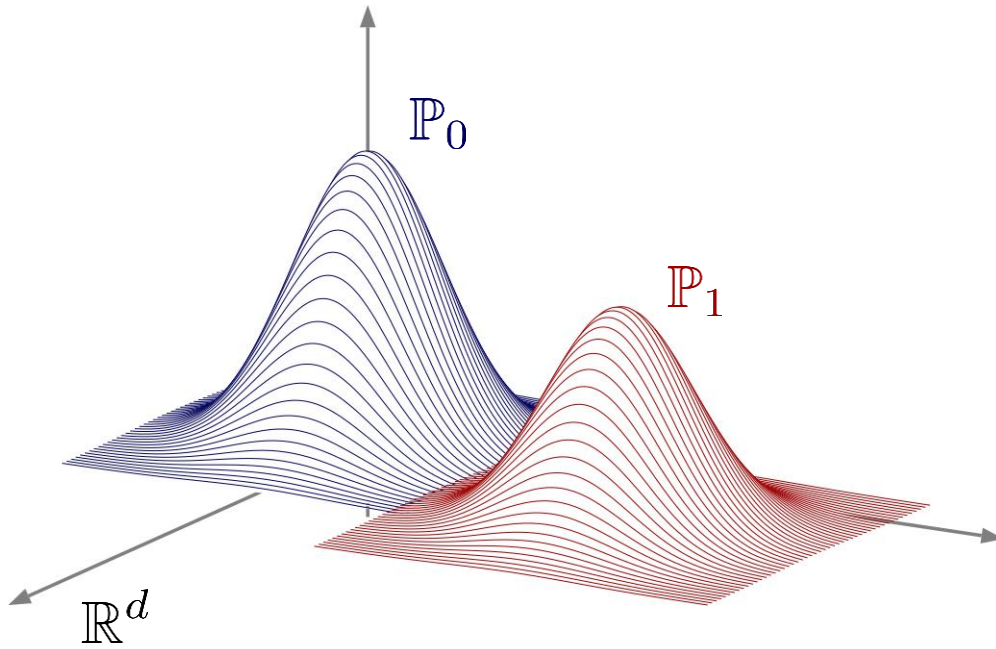


Problem

Observe a random vector $X \in \mathbb{R}^d$

Motivation

Hypothesis testing in high-dimensions



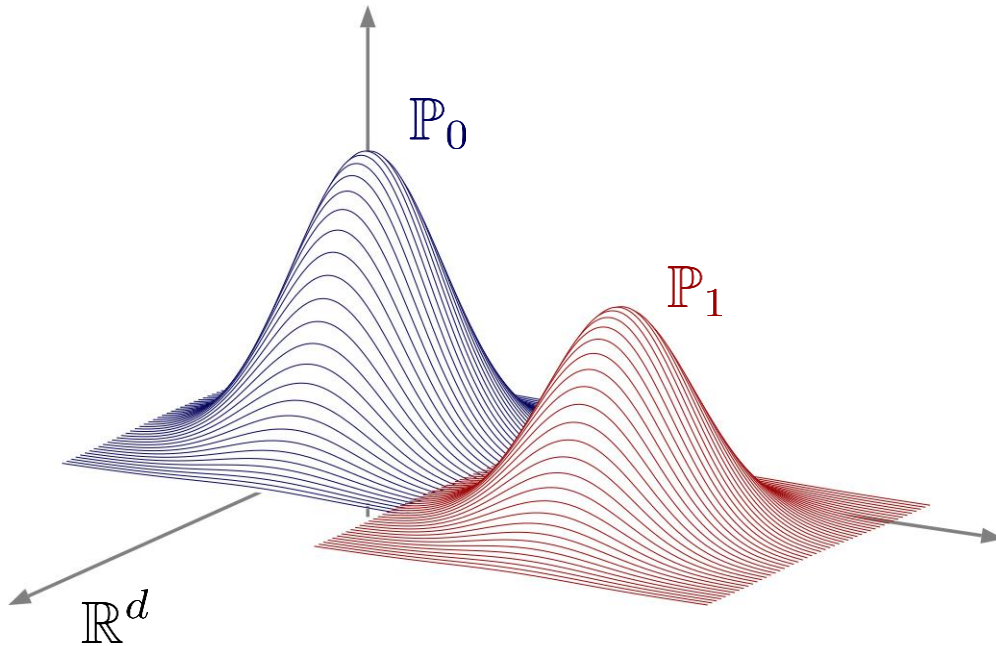
Problem

Observe a random vector $X \in \mathbb{R}^d$

$X \sim P_0$ or $X \sim P_1$?

Motivation

Hypothesis testing in high-dimensions



Problem

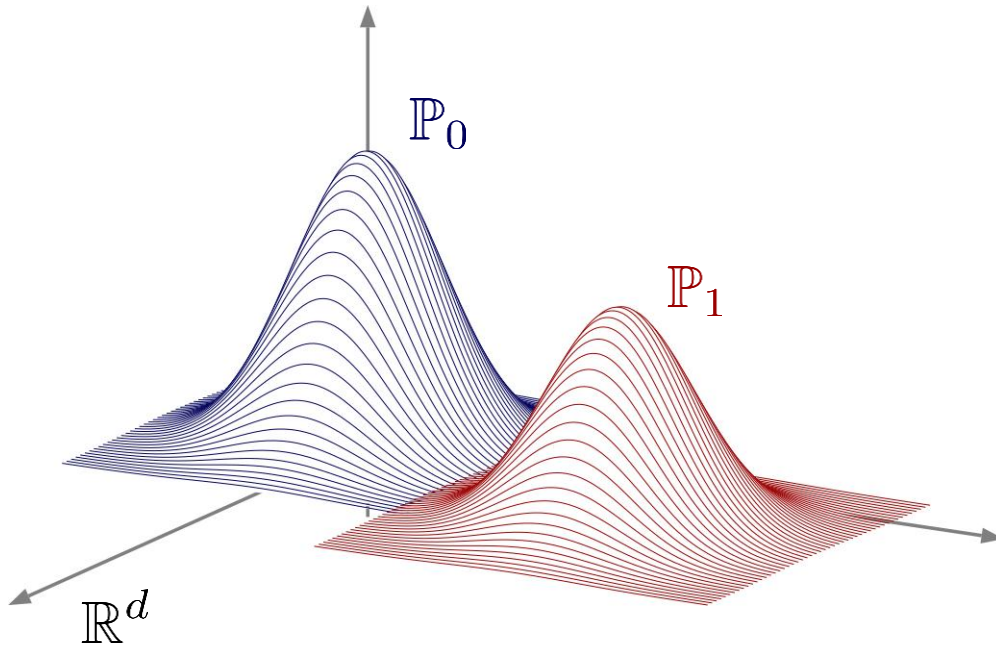
Observe a random vector $X \in \mathbb{R}^d$

$X \sim \mathbb{P}_0$ or $X \sim \mathbb{P}_1$?

False positive: $\alpha = \mathbb{P}(\text{decide } X \sim \mathbb{P}_1 \mid X \sim \mathbb{P}_0)$

Motivation

Hypothesis testing in high-dimensions



Problem

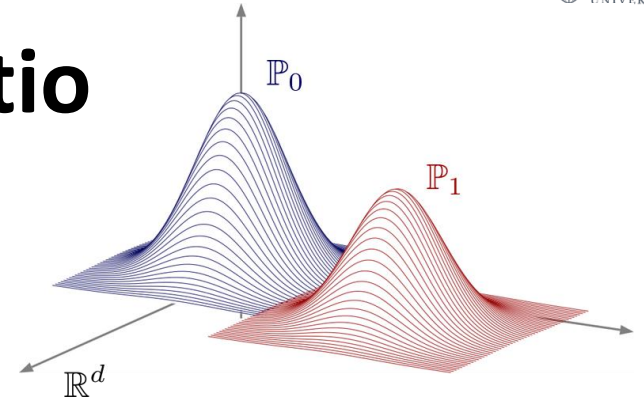
Observe a random vector $X \in \mathbb{R}^d$

$X \sim \mathbb{P}_0$ or $X \sim \mathbb{P}_1$?

False positive: $\alpha = \mathbb{P}(\text{decide } X \sim \mathbb{P}_1 \mid X \sim \mathbb{P}_0)$

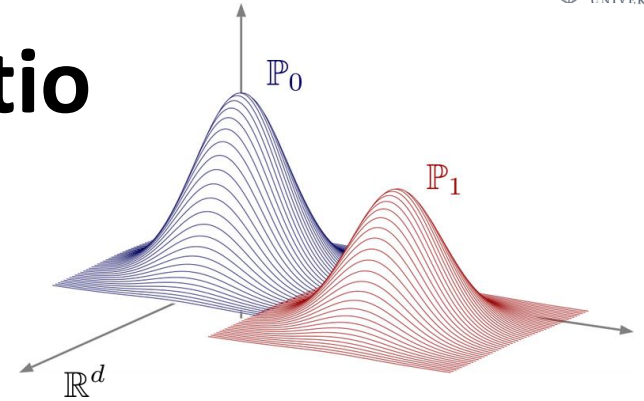
False negative: $\beta = \mathbb{P}(\text{decide } X \sim \mathbb{P}_0 \mid X \sim \mathbb{P}_1)$

Decision Rule: Likelihood Ratio



Decision Rule: Likelihood Ratio

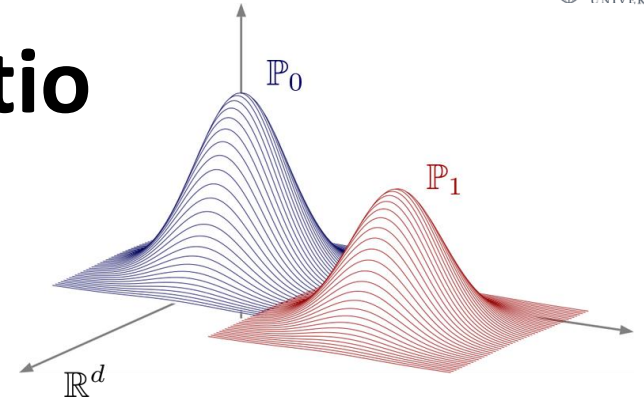
x_1, \dots, x_n : i.i.d. realizations of X



Decision Rule: Likelihood Ratio

x_1, \dots, x_n : i.i.d. realizations of X

For a given $T \geq 0$,

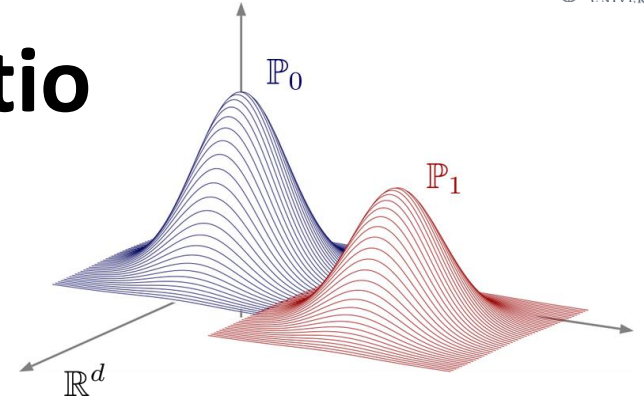


Decision Rule: Likelihood Ratio

x_1, \dots, x_n : i.i.d. realizations of X

For a given $T \geq 0$,

$$\text{If } \frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} > T$$

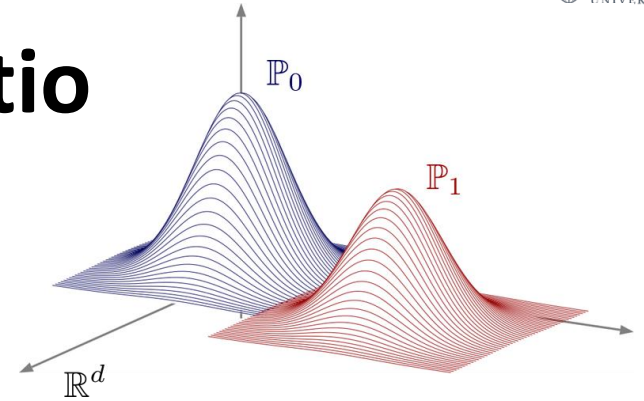


Decision Rule: Likelihood Ratio

x_1, \dots, x_n : i.i.d. realizations of X

For a given $T \geq 0$,

If $\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} > T$, then decide $X \sim \mathbb{P}_1$



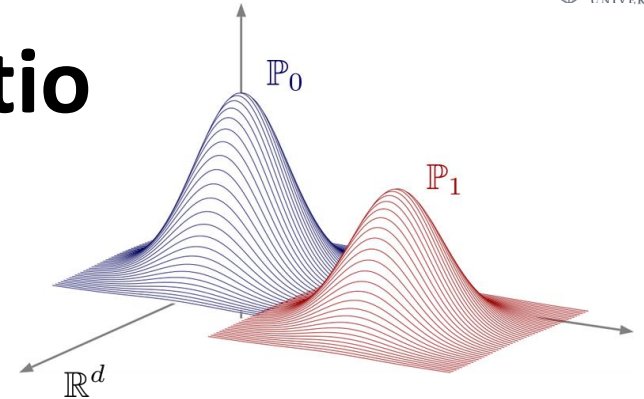
Decision Rule: Likelihood Ratio

x_1, \dots, x_n : i.i.d. realizations of X

For a given $T \geq 0$,

If $\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} > T$, then decide $X \sim \mathbb{P}_1$

If $\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} \leq T$, then decide $X \sim \mathbb{P}_0$



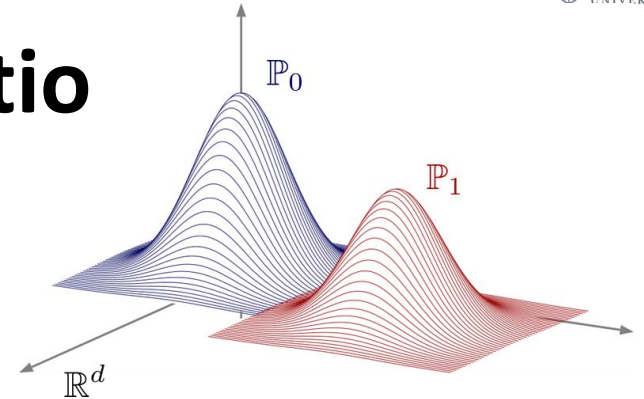
Decision Rule: Likelihood Ratio

x_1, \dots, x_n : i.i.d. realizations of X

For a given $T \geq 0$,

If $\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} > T$, then decide $X \sim \mathbb{P}_1$

If $\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} \leq T$, then decide $X \sim \mathbb{P}_0$



False positive: $\alpha_L = \mathbb{P}(L(x_1, \dots, x_n) > T \mid X \sim \mathbb{P}_0)$

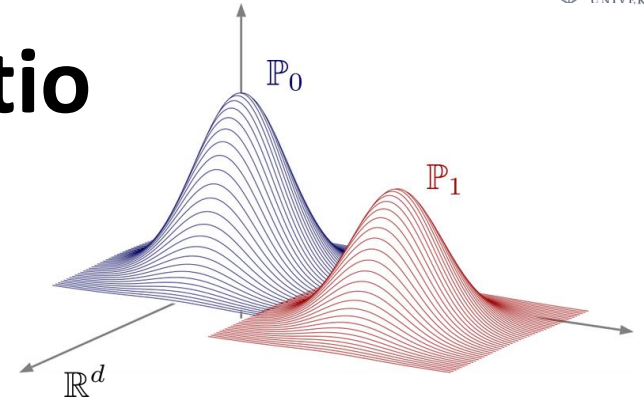
Decision Rule: Likelihood Ratio

x_1, \dots, x_n : i.i.d. realizations of X

For a given $T \geq 0$,

If $\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} > T$, then decide $X \sim \mathbb{P}_1$

If $\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} \leq T$, then decide $X \sim \mathbb{P}_0$



$$\text{False positive: } \alpha_L = \mathbb{P}\left(\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} > T \mid X \sim \mathbb{P}_0 \right)$$

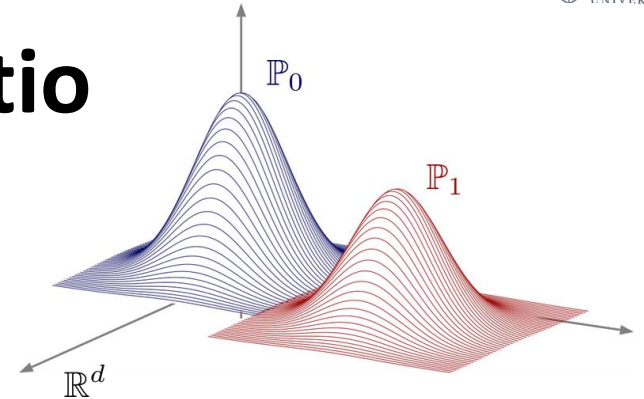
Decision Rule: Likelihood Ratio

x_1, \dots, x_n : i.i.d. realizations of X

For a given $T \geq 0$,

If $\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} > T$, then decide $X \sim \mathbb{P}_1$

If $\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} \leq T$, then decide $X \sim \mathbb{P}_0$



False positive: $\alpha_L = \mathbb{P}\left(\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} > T \mid X \sim \mathbb{P}_0 \right)$

False negative: $\beta_L = \mathbb{P}\left(\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} \leq T \mid X \sim \mathbb{P}_1 \right)$

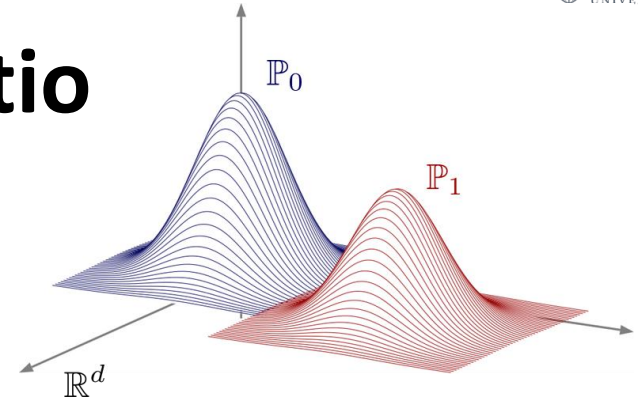
Decision Rule: Likelihood Ratio

x_1, \dots, x_n : i.i.d. realizations of X

For a given $T \geq 0$,

If $\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} > T$, then decide $X \sim \mathbb{P}_1$

If $\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)} \leq T$, then decide $X \sim \mathbb{P}_0$



MAP rule / minimizes risk when

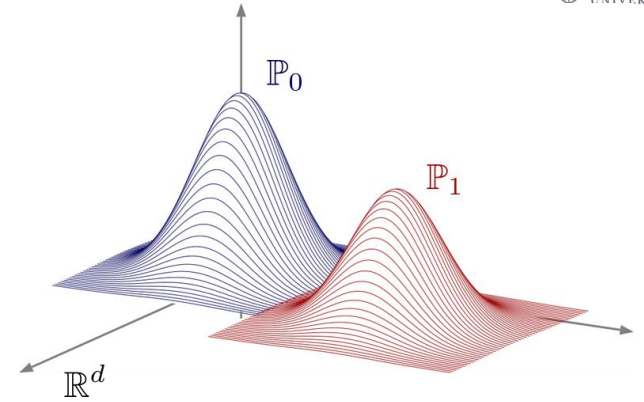
$$T = \frac{\mathbb{P}(X \sim \mathbb{P}_0)}{\mathbb{P}(X \sim \mathbb{P}_1)}$$

$$\frac{\mathbb{P}_1(x_1, \dots, x_n)}{\mathbb{P}_0(x_1, \dots, x_n)}$$

False positive: $\alpha_L = \mathbb{P}(L(x_1, \dots, x_n) > T \mid X \sim \mathbb{P}_0)$

False negative: $\beta_L = \mathbb{P}(L(x_1, \dots, x_n) \leq T \mid X \sim \mathbb{P}_1)$

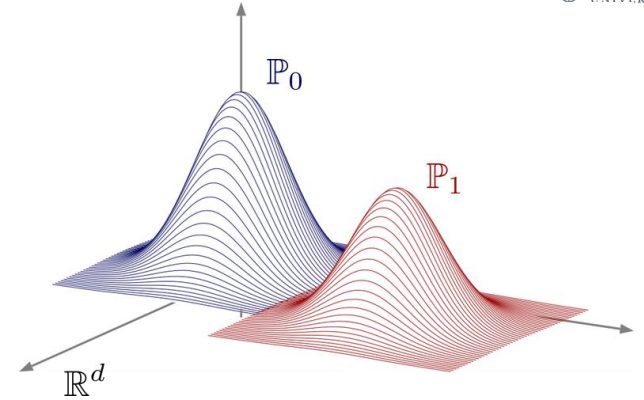
Neyman-Pearson Lemma



Neyman-Pearson Lemma

$$\alpha_L = \mathbb{P}(L(x_1, \dots, x_n) > T \mid X \sim \mathbb{P}_0)$$

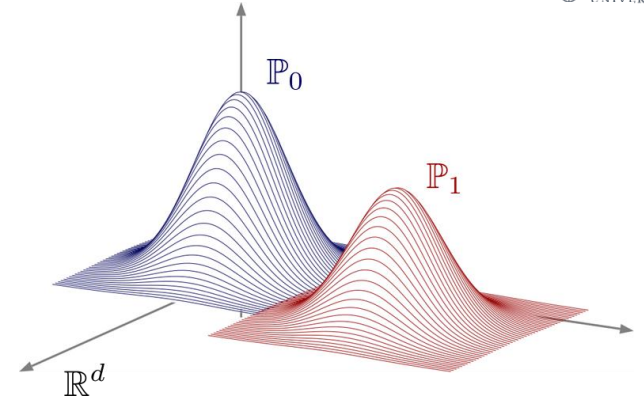
$$\beta_L = \mathbb{P}(L(x_1, \dots, x_n) \leq T \mid X \sim \mathbb{P}_1)$$



Neyman-Pearson Lemma

$$\alpha_L = \mathbb{P}(L(x_1, \dots, x_n) > T \mid X \sim \mathbb{P}_0)$$

$$\beta_L = \mathbb{P}(L(x_1, \dots, x_n) \leq T \mid X \sim \mathbb{P}_1)$$

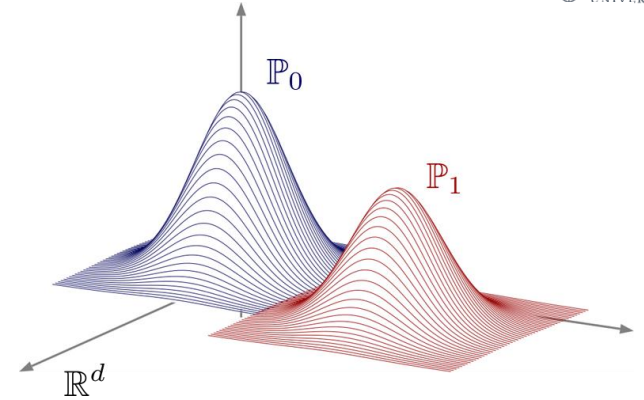


Neyman-Pearson Lemma

Neyman-Pearson Lemma

$$\alpha_L = \mathbb{P}(L(x_1, \dots, x_n) > T \mid X \sim \mathbb{P}_0)$$

$$\beta_L = \mathbb{P}(L(x_1, \dots, x_n) \leq T \mid X \sim \mathbb{P}_1)$$



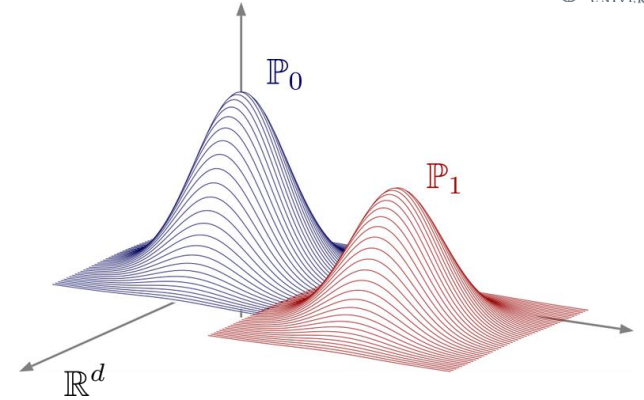
Neyman-Pearson Lemma

The likelihood ratio test is *optimal*:

Neyman-Pearson Lemma

$$\alpha_L = \mathbb{P}(L(x_1, \dots, x_n) > T \mid X \sim \mathbb{P}_0)$$

$$\beta_L = \mathbb{P}(L(x_1, \dots, x_n) \leq T \mid X \sim \mathbb{P}_1)$$



Neyman-Pearson Lemma

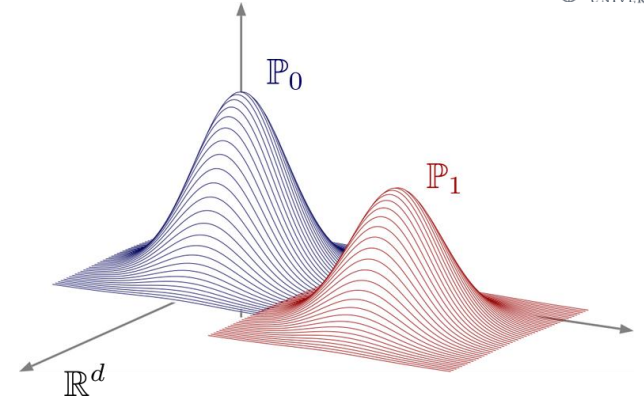
The likelihood ratio test is *optimal*:

If there is another (possibly random) decision rule $D(x_1, \dots, x_n)$ such that

Neyman-Pearson Lemma

$$\alpha_L = \mathbb{P}(L(x_1, \dots, x_n) > T \mid X \sim \mathbb{P}_0)$$

$$\beta_L = \mathbb{P}(L(x_1, \dots, x_n) \leq T \mid X \sim \mathbb{P}_1)$$



Neyman-Pearson Lemma

The likelihood ratio test is *optimal*:

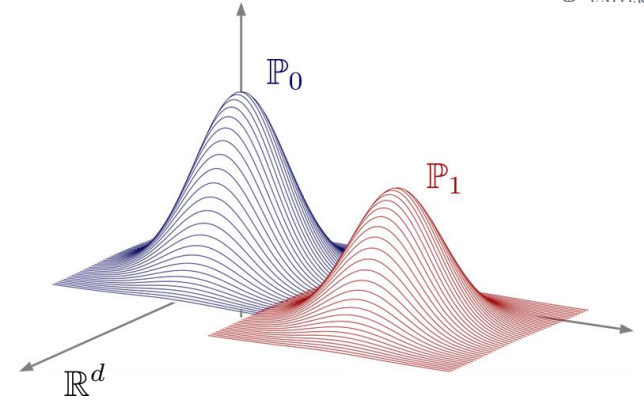
If there is another (possibly random) decision rule $D(x_1, \dots, x_n)$ such that

$$\mathbb{P}(D(x_1, \dots, x_n) \text{ decides } X \sim \mathbb{P}_1 \mid X \sim \mathbb{P}_0) \leq \alpha_L,$$

Neyman-Pearson Lemma

$$\alpha_L = \mathbb{P}(L(x_1, \dots, x_n) > T \mid X \sim \mathbb{P}_0)$$

$$\beta_L = \mathbb{P}(L(x_1, \dots, x_n) \leq T \mid X \sim \mathbb{P}_1)$$



Neyman-Pearson Lemma

The likelihood ratio test is *optimal*:

If there is another (possibly random) decision rule $D(x_1, \dots, x_n)$ such that

$$\mathbb{P}(D(x_1, \dots, x_n) \text{ decides } X \sim \mathbb{P}_1 \mid X \sim \mathbb{P}_0) \leq \alpha_L,$$

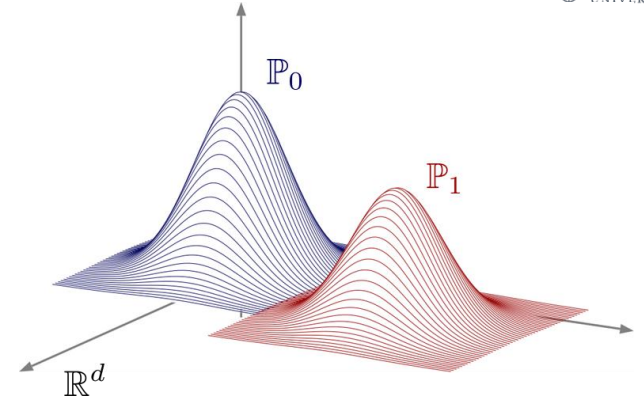
then

$$\mathbb{P}(D(x_1, \dots, x_n) \text{ decides } X \sim \mathbb{P}_0 \mid X \sim \mathbb{P}_1) \geq \beta_L.$$

Neyman-Pearson Lemma

$$\alpha_L = \mathbb{P}(L(x_1, \dots, x_n) > T \mid X \sim \mathbb{P}_0)$$

$$\beta_L = \mathbb{P}(L(x_1, \dots, x_n) \leq T \mid X \sim \mathbb{P}_1)$$



Neyman-Pearson Lemma

The likelihood ratio test is *optimal*:

If there is another (possibly random) decision rule $D(x_1, \dots, x_n)$ such that

$$\mathbb{P}(D(x_1, \dots, x_n) \text{ decides } X \sim \mathbb{P}_1 \mid X \sim \mathbb{P}_0) \leq \alpha_L,$$

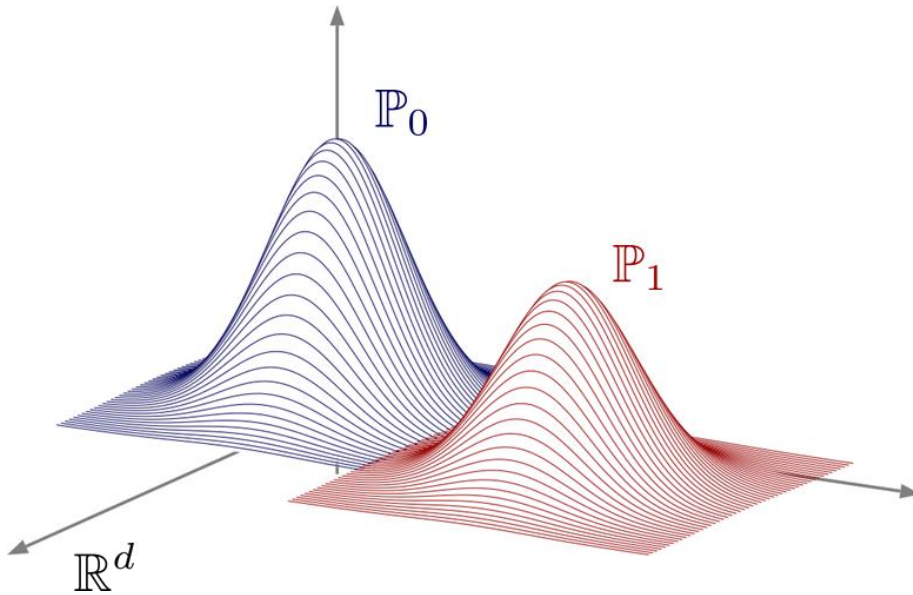
then

$$\mathbb{P}(D(x_1, \dots, x_n) \text{ decides } X \sim \mathbb{P}_0 \mid X \sim \mathbb{P}_1) \geq \beta_L.$$

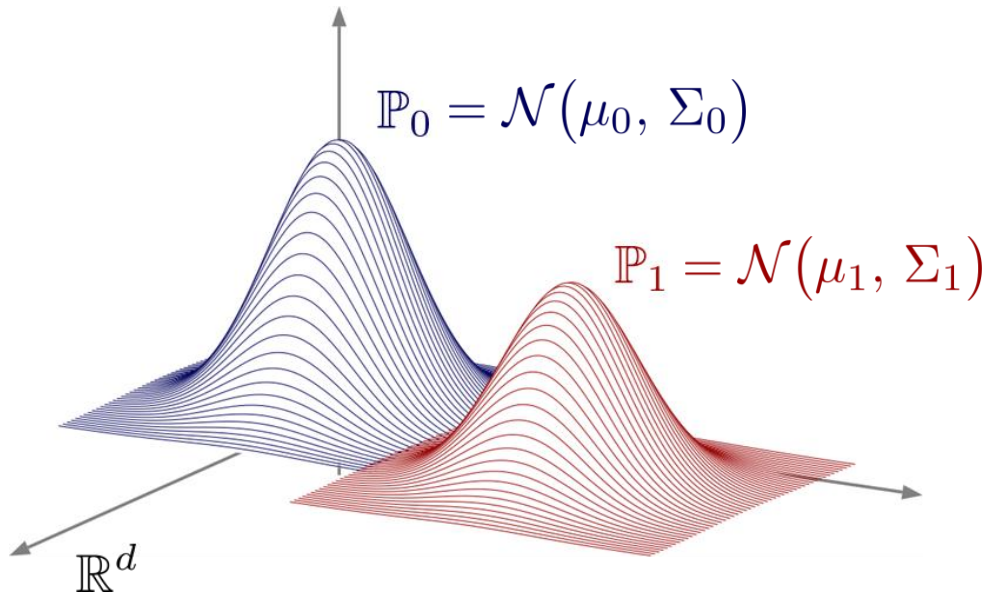
And vice-versa.

Linear Discriminant Analysis

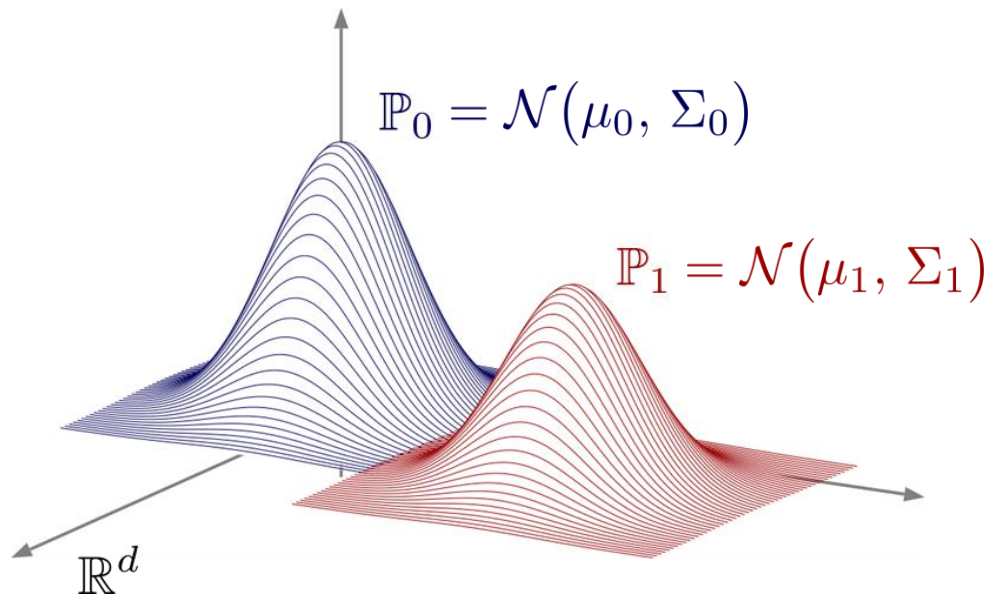
Linear Discriminant Analysis



Linear Discriminant Analysis

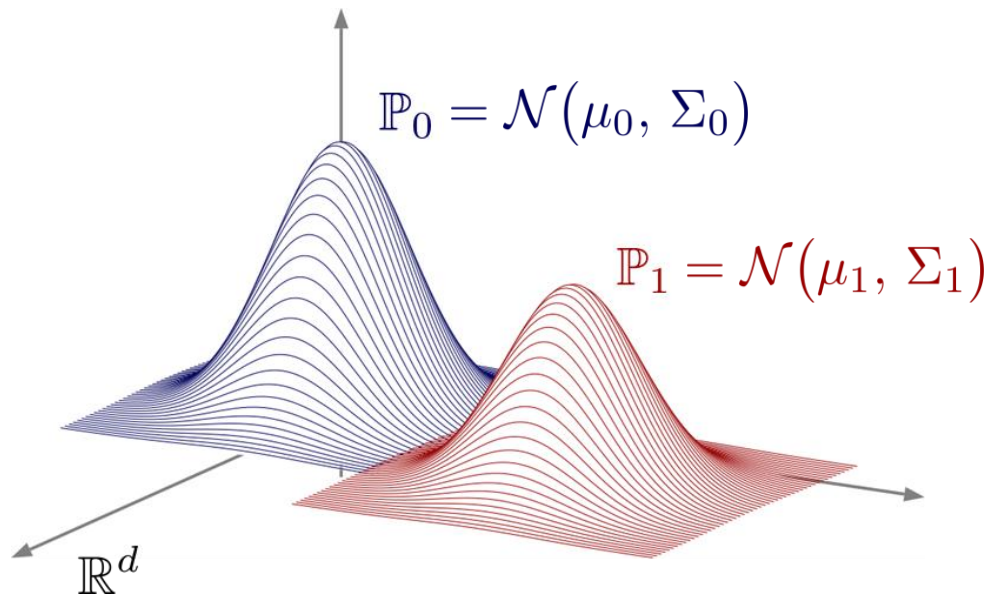


Linear Discriminant Analysis



$$X \sim \mathcal{N}(\mu, \Sigma) \quad \Rightarrow \quad f_X(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right)$$

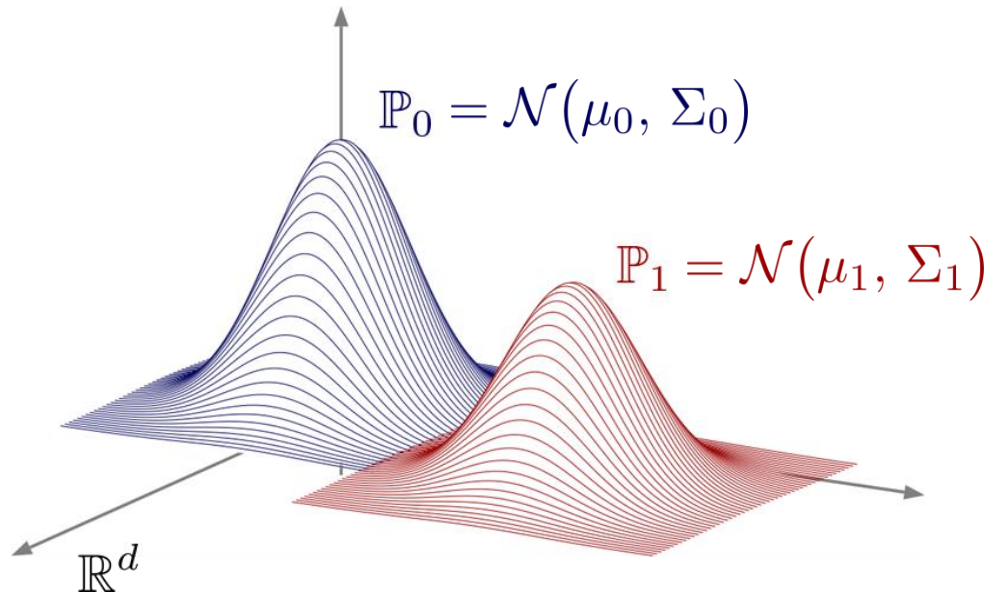
Linear Discriminant Analysis



$$X \sim \mathcal{N}(\mu, \Sigma) \quad \Rightarrow \quad f_X(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right)$$

Assume (to simplify): $\Sigma_0 = \Sigma_1 = \Sigma$

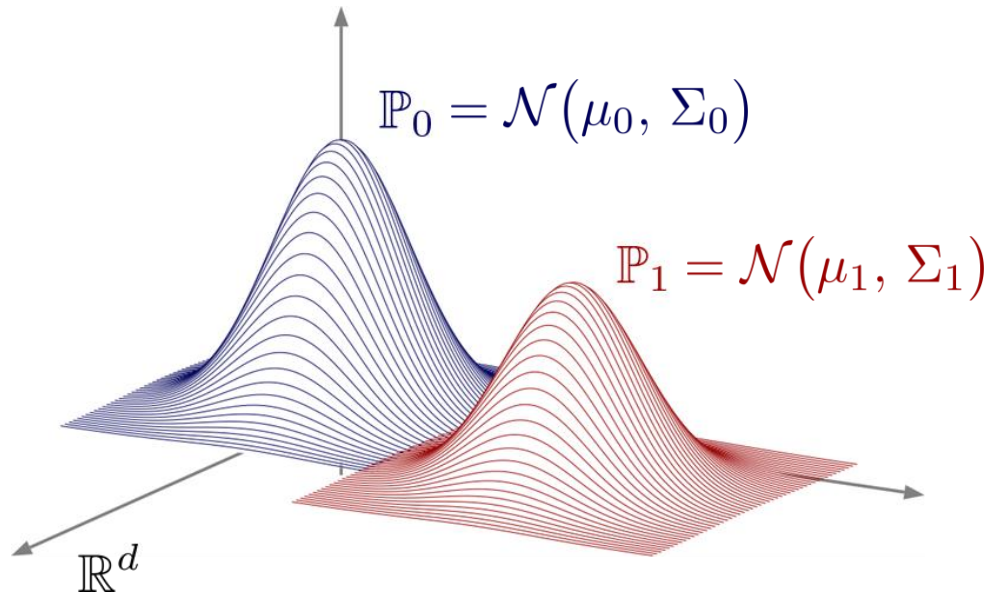
Linear Discriminant Analysis



$$X \sim \mathcal{N}(\mu, \Sigma) \quad \Rightarrow \quad f_X(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right)$$

Assume (to simplify): $\Sigma_0 = \Sigma_1 = \Sigma$ and $n = 1$ (one observation)

Linear Discriminant Analysis

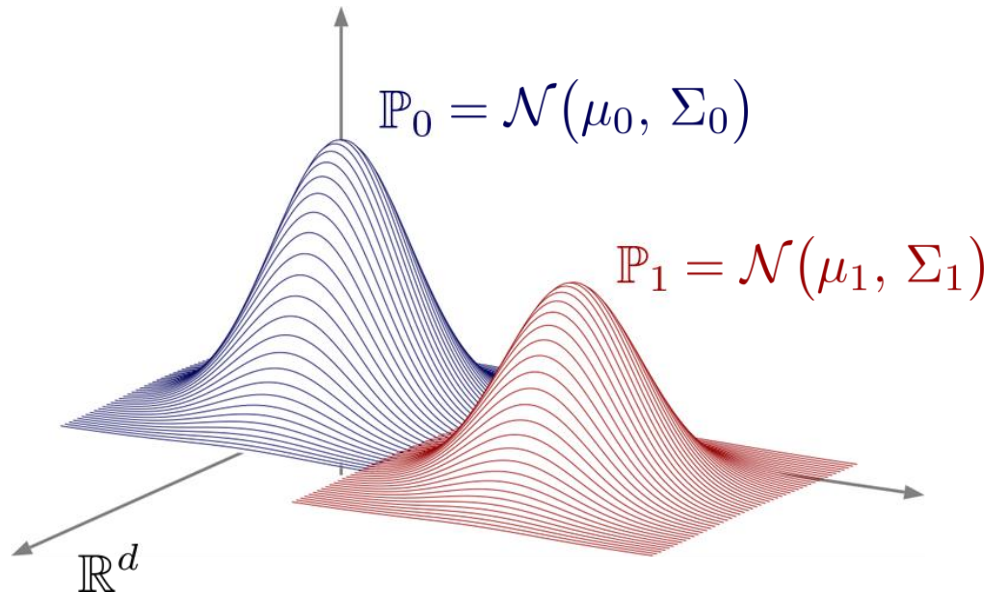


$$X \sim \mathcal{N}(\mu, \Sigma) \quad \Rightarrow \quad f_X(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right)$$

Assume (to simplify): $\Sigma_0 = \Sigma_1 = \Sigma$ and $n = 1$ (one observation)

$$L(x) > T$$

Linear Discriminant Analysis



$$X \sim \mathcal{N}(\mu, \Sigma) \quad \Rightarrow \quad f_X(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right)$$

Assume (to simplify): $\Sigma_0 = \Sigma_1 = \Sigma$ and $n = 1$ (one observation)

$$L(x) > T \quad \Longleftrightarrow \quad \Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_0) > \log T$$

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2}\right)^\top \Sigma^{-1}(\mu_1 - \mu_0) > \log T$$

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2}\right)^\top \Sigma^{-1}(\mu_1 - \mu_0) > \log T$$

Probability of error

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_0) > \log T$$

Probability of error (assuming \mathbb{P}_0 and \mathbb{P}_1 are equally likely)

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_0) > \log T$$

Probability of error (assuming \mathbb{P}_0 and \mathbb{P}_1 are equally likely)

$$\text{Err}(\Psi) = \mathbb{P}(\Psi(X) > 0 \ \& \ \mathbb{P}_0 \text{ true}) + \mathbb{P}(\Psi(X) \leq 0 \ \& \ \mathbb{P}_1 \text{ true})$$

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_0) > \log T$$

Probability of error (assuming \mathbb{P}_0 and \mathbb{P}_1 are equally likely)

$$\begin{aligned} \text{Err}(\Psi) &= \mathbb{P}(\Psi(X) > 0 \ \& \ \mathbb{P}_0 \text{ true}) + \mathbb{P}(\Psi(X) \leq 0 \ \& \ \mathbb{P}_1 \text{ true}) \\ &= \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) \cdot \mathbb{P}(\mathbb{P}_0) + \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \cdot \mathbb{P}(\mathbb{P}_1) \end{aligned}$$

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_0) > \log T$$

Probability of error (assuming \mathbb{P}_0 and \mathbb{P}_1 are equally likely)

$$\begin{aligned} \text{Err}(\Psi) &= \mathbb{P}(\Psi(X) > 0 \ \& \ \mathbb{P}_0 \text{ true}) + \mathbb{P}(\Psi(X) \leq 0 \ \& \ \mathbb{P}_1 \text{ true}) \\ &= \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) \cdot \mathbb{P}(\mathbb{P}_0) + \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \cdot \mathbb{P}(\mathbb{P}_1) \\ &= \frac{1}{2} \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) + \frac{1}{2} \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \end{aligned}$$

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_0) > \log T$$

Probability of error (assuming \mathbb{P}_0 and \mathbb{P}_1 are equally likely)

$$\begin{aligned} \text{Err}(\Psi) &= \mathbb{P}(\Psi(X) > 0 \ \& \ \mathbb{P}_0 \text{ true}) + \mathbb{P}(\Psi(X) \leq 0 \ \& \ \mathbb{P}_1 \text{ true}) \\ &= \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) \cdot \mathbb{P}(\mathbb{P}_0) + \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \cdot \mathbb{P}(\mathbb{P}_1) \\ &= \frac{1}{2} \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) + \frac{1}{2} \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \end{aligned}$$

(using Gaussianity and manipulating...)

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_0) > \log T$$

Probability of error (assuming \mathbb{P}_0 and \mathbb{P}_1 are equally likely)

$$\begin{aligned} \text{Err}(\Psi) &= \mathbb{P}(\Psi(X) > 0 \ \& \ \mathbb{P}_0 \text{ true}) + \mathbb{P}(\Psi(X) \leq 0 \ \& \ \mathbb{P}_1 \text{ true}) \\ &= \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) \cdot \mathbb{P}(\mathbb{P}_0) + \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \cdot \mathbb{P}(\mathbb{P}_1) \\ &= \frac{1}{2} \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) + \frac{1}{2} \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \end{aligned}$$

(using Gaussianity and manipulating...)

$$= \Phi\left(-\frac{\gamma}{2}\right) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-\gamma/2} e^{-t^2/2} dt$$

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_0) > \log T$$

Probability of error (assuming \mathbb{P}_0 and \mathbb{P}_1 are equally likely)

$$\begin{aligned} \text{Err}(\Psi) &= \mathbb{P}(\Psi(X) > 0 \ \& \ \mathbb{P}_0 \text{ true}) + \mathbb{P}(\Psi(X) \leq 0 \ \& \ \mathbb{P}_1 \text{ true}) \\ &= \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) \cdot \mathbb{P}(\mathbb{P}_0) + \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \cdot \mathbb{P}(\mathbb{P}_1) \\ &= \frac{1}{2} \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) + \frac{1}{2} \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \end{aligned}$$

(using Gaussianity and manipulating...)

$$= \Phi\left(-\frac{\gamma}{2}\right) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-\gamma/2} e^{-t^2/2} dt \qquad \gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1)}$$

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2}\right)^\top \Sigma^{-1}(\mu_1 - \mu_0) > \log T$$

Probability of error (assuming \mathbb{P}_0 and \mathbb{P}_1 are equally likely)

$$\begin{aligned}\text{Err}(\Psi) &= \mathbb{P}(\Psi(X) > 0 \ \& \ \mathbb{P}_0 \text{ true}) + \mathbb{P}(\Psi(X) \leq 0 \ \& \ \mathbb{P}_1 \text{ true}) \\ &= \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) \cdot \mathbb{P}(\mathbb{P}_0) + \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \cdot \mathbb{P}(\mathbb{P}_1) \\ &= \frac{1}{2} \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) + \frac{1}{2} \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1)\end{aligned}$$

(using Gaussianity and manipulating...)

$$= \Phi\left(-\frac{\gamma}{2}\right) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-\gamma/2} e^{-t^2/2} dt$$

classical error expression

$$\gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1}(\mu_0 - \mu_1)}$$

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_0) > \log T$$

Probability of error (assuming \mathbb{P}_0 and \mathbb{P}_1 are equally likely)

$$\begin{aligned} \text{Err}(\Psi) &= \mathbb{P}(\Psi(X) > 0 \ \& \ \mathbb{P}_0 \text{ true}) + \mathbb{P}(\Psi(X) \leq 0 \ \& \ \mathbb{P}_1 \text{ true}) \\ &= \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) \cdot \mathbb{P}(\mathbb{P}_0) + \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \cdot \mathbb{P}(\mathbb{P}_1) \\ &= \frac{1}{2} \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) + \frac{1}{2} \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \end{aligned}$$

(using Gaussianity and manipulating...)

$$= \Phi\left(-\frac{\gamma}{2}\right) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-\gamma/2} e^{-t^2/2} dt$$

classical error expression

$$\gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1)}$$

$\begin{array}{ccc} | & | & | \\ \mu_0 & - & \mu_1 \end{array}$
 need to be estimated

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_0) > \log T$$

Probability of error (assuming \mathbb{P}_0 and \mathbb{P}_1 are equally likely)

$$\begin{aligned} \text{Err}(\Psi) &= \mathbb{P}(\Psi(X) > 0 \ \& \ \mathbb{P}_0 \text{ true}) + \mathbb{P}(\Psi(X) \leq 0 \ \& \ \mathbb{P}_1 \text{ true}) \\ &= \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) \cdot \mathbb{P}(\mathbb{P}_0) + \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \cdot \mathbb{P}(\mathbb{P}_1) \\ &= \frac{1}{2} \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) + \frac{1}{2} \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \end{aligned}$$

(using Gaussianity and manipulating...)

$$= \Phi\left(-\frac{\gamma}{2}\right) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-\gamma/2} e^{-t^2/2} dt$$

classical error expression

$$\gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1)}$$

$\begin{array}{c} | \quad | \quad | \\ \text{need to be estimated} \\ n_0 \text{ and } n_1 \text{ samples} \end{array}$

$$\Psi(x) := \left(x - \frac{\mu_0 + \mu_1}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_0) > \log T$$

Probability of error (assuming \mathbb{P}_0 and \mathbb{P}_1 are equally likely)

$$\begin{aligned} \text{Err}(\Psi) &= \mathbb{P}(\Psi(X) > 0 \ \& \ \mathbb{P}_0 \text{ true}) + \mathbb{P}(\Psi(X) \leq 0 \ \& \ \mathbb{P}_1 \text{ true}) \\ &= \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) \cdot \mathbb{P}(\mathbb{P}_0) + \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \cdot \mathbb{P}(\mathbb{P}_1) \\ &= \frac{1}{2} \mathbb{P}(\Psi(X) > 0 \mid \mathbb{P}_0) + \frac{1}{2} \mathbb{P}(\Psi(X) \leq 0 \mid \mathbb{P}_1) \end{aligned}$$

(using Gaussianity and manipulating...)

$$= \Phi\left(-\frac{\gamma}{2}\right) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-\gamma/2} e^{-t^2/2} dt$$

classical error expression

$$\gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1)}$$

need to be estimated

n_0 and n_1 samples

High-dimensional regime: n_0 and n_1 same order as d

Fisher Linear Discriminant

Fisher Linear Discriminant

$$\Psi(x) = \left(x - \frac{\mu_0 + \mu_1}{2}\right)^\top \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1}(\mu_0 - \mu_1)} \quad \text{Err}(\Psi) = \Phi\left(-\frac{\gamma}{2}\right)$$

Fisher Linear Discriminant

$$\Psi(x) = \left(x - \frac{\mu_0 + \mu_1}{2}\right)^\top \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1}(\mu_0 - \mu_1)} \quad \text{Err}(\Psi) = \Phi\left(-\frac{\gamma}{2}\right)$$

Unbiased estimators:

Fisher Linear Discriminant

$$\Psi(x) = \left(x - \frac{\mu_0 + \mu_1}{2}\right)^\top \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1}(\mu_0 - \mu_1)} \quad \text{Err}(\Psi) = \Phi\left(-\frac{\gamma}{2}\right)$$

Unbiased estimators:

$$\hat{\mu}_0 := \frac{1}{n_0} \sum_{i=1}^{n_0} x_i$$

$$\hat{\mu}_1 := \frac{1}{n_1} \sum_{i=1}^{n_1} y_i$$

Fisher Linear Discriminant

$$\Psi(x) = \left(x - \frac{\mu_0 + \mu_1}{2}\right)^\top \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1}(\mu_0 - \mu_1)} \quad \text{Err}(\Psi) = \Phi\left(-\frac{\gamma}{2}\right)$$

Unbiased estimators:

$$\hat{\mu}_0 := \frac{1}{n_0} \sum_{i=1}^{n_0} x_i$$

$$\hat{\mu}_1 := \frac{1}{n_1} \sum_{i=1}^{n_1} y_i$$

$$\hat{\Sigma} := \frac{1}{n_0 - 1} \sum_{i=1}^{n_0} (x_i - \hat{\mu}_0)(x_i - \hat{\mu}_0)^\top + \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (y_i - \hat{\mu}_1)(y_i - \hat{\mu}_1)^\top$$

Fisher Linear Discriminant

$$\Psi(x) = \left(x - \frac{\mu_0 + \mu_1}{2}\right)^\top \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1}(\mu_0 - \mu_1)} \quad \text{Err}(\Psi) = \Phi\left(-\frac{\gamma}{2}\right)$$

Unbiased estimators:

$$\hat{\mu}_0 := \frac{1}{n_0} \sum_{i=1}^{n_0} x_i$$

$$\hat{\mu}_1 := \frac{1}{n_1} \sum_{i=1}^{n_1} y_i$$

$$\hat{\Sigma} := \frac{1}{n_0 - 1} \sum_{i=1}^{n_0} (x_i - \hat{\mu}_0)(x_i - \hat{\mu}_0)^\top + \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (y_i - \hat{\mu}_1)(y_i - \hat{\mu}_1)^\top$$

Plug estimators into log-likelihood ratio:

Fisher Linear Discriminant

$$\Psi(x) = \left(x - \frac{\mu_0 + \mu_1}{2}\right)^\top \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1}(\mu_0 - \mu_1)} \quad \text{Err}(\Psi) = \Phi\left(-\frac{\gamma}{2}\right)$$

Unbiased estimators:

$$\hat{\mu}_0 := \frac{1}{n_0} \sum_{i=1}^{n_0} x_i$$

$$\hat{\mu}_1 := \frac{1}{n_1} \sum_{i=1}^{n_1} y_i$$

$$\hat{\Sigma} := \frac{1}{n_0 - 1} \sum_{i=1}^{n_0} (x_i - \hat{\mu}_0)(x_i - \hat{\mu}_0)^\top + \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (y_i - \hat{\mu}_1)(y_i - \hat{\mu}_1)^\top$$

Plug estimators into log-likelihood ratio:

$$\hat{\Psi}(x) := \left(x - \frac{\hat{\mu}_0 + \hat{\mu}_1}{2}\right)^\top \hat{\Sigma}^{-1}(\hat{\mu}_1 - \hat{\mu}_0)$$

Fisher Linear Discriminant

$$\Psi(x) = \left(x - \frac{\mu_0 + \mu_1}{2}\right)^\top \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\gamma = \sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1}(\mu_0 - \mu_1)} \quad \text{Err}(\Psi) = \Phi\left(-\frac{\gamma}{2}\right)$$

Unbiased estimators:

$$\hat{\mu}_0 := \frac{1}{n_0} \sum_{i=1}^{n_0} x_i$$

$$\hat{\mu}_1 := \frac{1}{n_1} \sum_{i=1}^{n_1} y_i$$

$$\hat{\Sigma} := \frac{1}{n_0 - 1} \sum_{i=1}^{n_0} (x_i - \hat{\mu}_0)(x_i - \hat{\mu}_0)^\top + \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (y_i - \hat{\mu}_1)(y_i - \hat{\mu}_1)^\top$$

Plug estimators into log-likelihood ratio:

$$\hat{\Psi}(x) := \left(x - \frac{\hat{\mu}_0 + \hat{\mu}_1}{2}\right)^\top \hat{\Sigma}^{-1}(\hat{\mu}_1 - \hat{\mu}_0)$$

Fisher linear discriminant function

Fisher Linear Discriminant

Fisher Linear Discriminant

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$ $d = 400$ $n_0 = n_1 = 800$

Fisher Linear Discriminant

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$ $d = 400$ $n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Fisher Linear Discriminant

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$ $d = 400$ $n_0 = n_1 = 800$

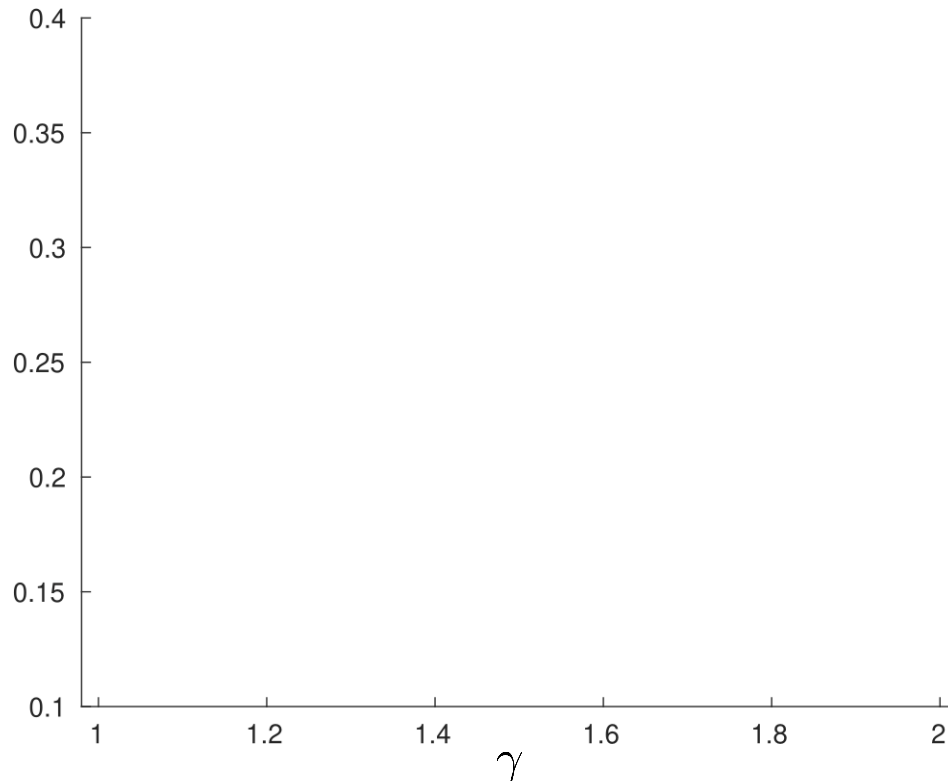
Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2 Test with $\hat{\Psi}(x)$ over 5000 random trials

Fisher Linear Discriminant

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$ $d = 400$ $n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2 Test with $\hat{\Psi}(x)$ over 5000 random trials

Probability of error



Fisher Linear Discriminant

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

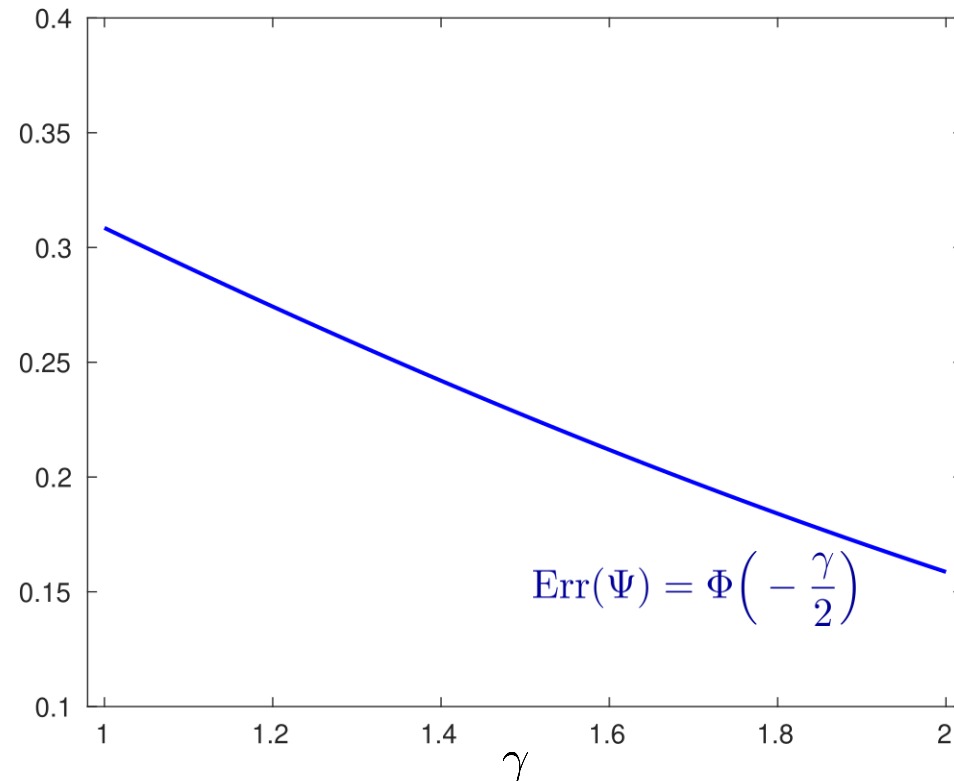
$d = 400$

$n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\hat{\Psi}(x)$ over 5000 random trials

Probability of error



Fisher Linear Discriminant

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

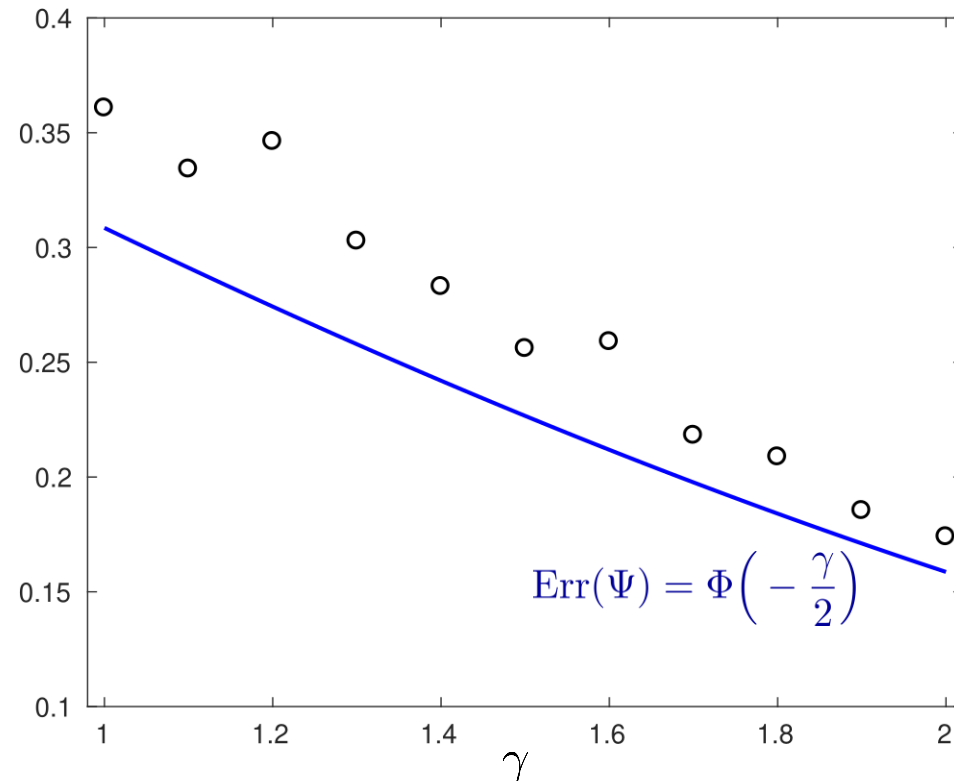
$d = 400$

$n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\hat{\Psi}(x)$ over 5000 random trials

Probability of error



Fisher Linear Discriminant

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

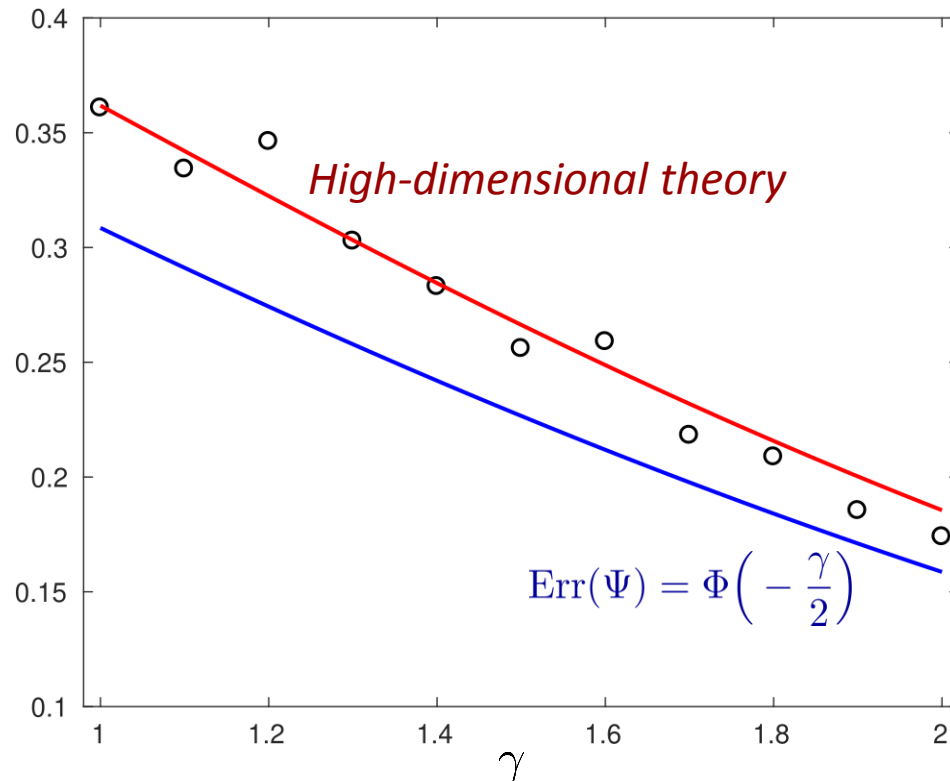
$d = 400$

$n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\hat{\Psi}(x)$ over 5000 random trials

Probability of error



Fisher Linear Discriminant

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

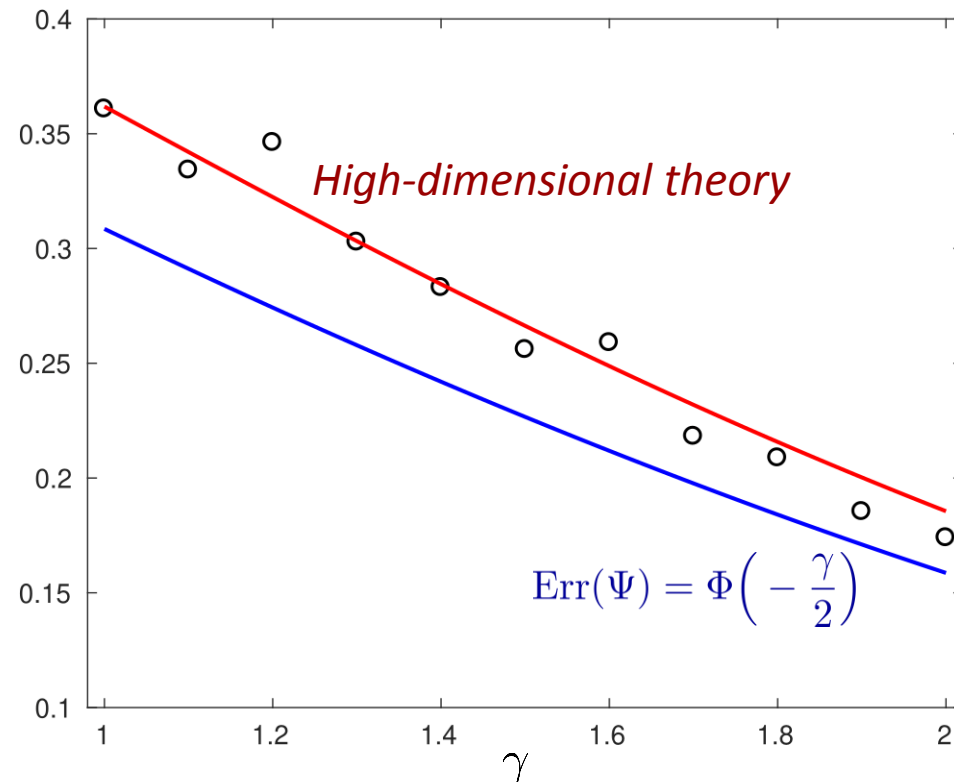
$d = 400$

$n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\hat{\Psi}(x)$ over 5000 random trials

Probability of error



[Kolmogorov]

Fisher Linear Discriminant

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

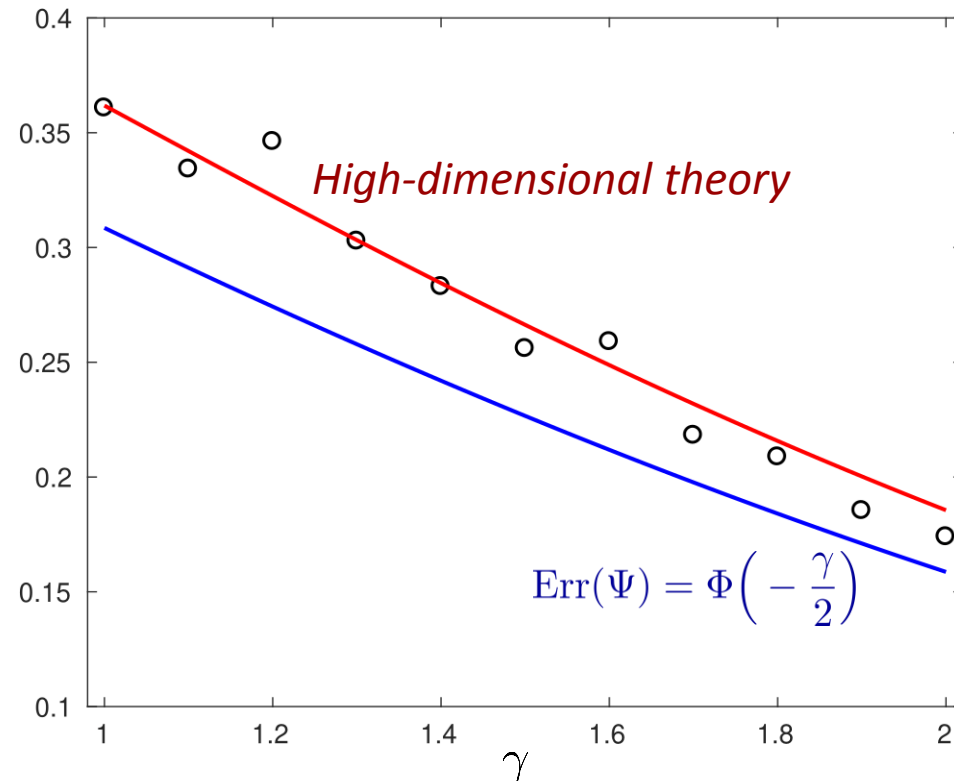
$d = 400$

$n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\hat{\Psi}(x)$ over 5000 random trials

Probability of error



[Kolmogorov]

$$(d, n_0, n_1) \rightarrow \infty \quad \frac{d}{n_0}, \frac{d}{n_1} \rightarrow \alpha$$

Fisher Linear Discriminant

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

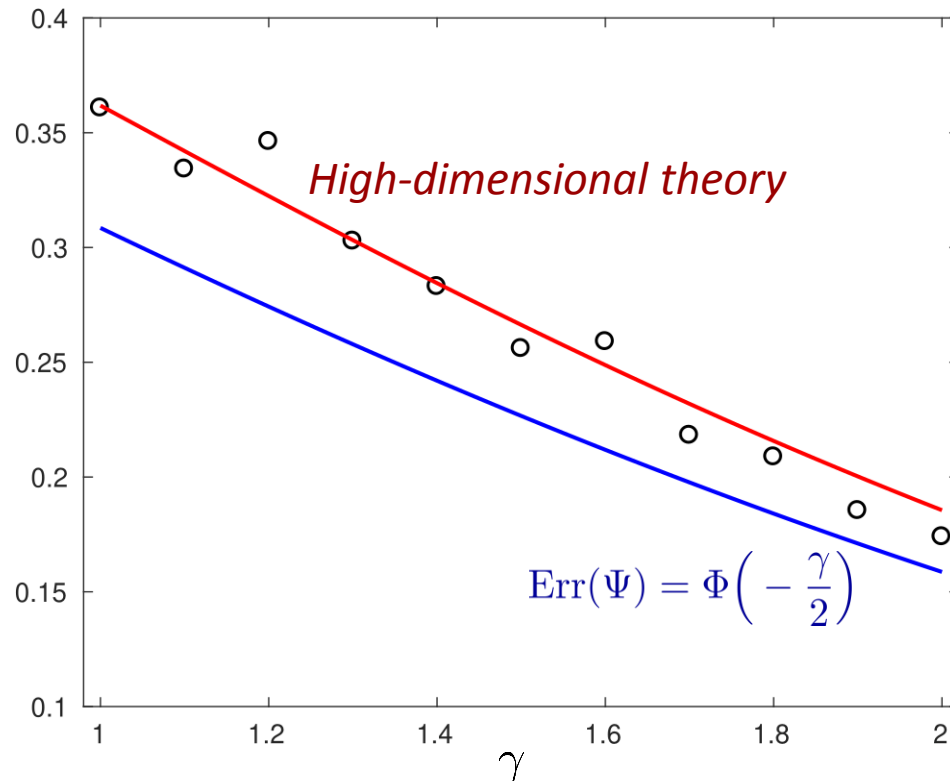
$d = 400$

$n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\widehat{\Psi}(x)$ over 5000 random trials

Probability of error



[Kolmogorov]

$$(d, n_0, n_1) \rightarrow \infty \quad \frac{d}{n_0}, \frac{d}{n_1} \rightarrow \alpha$$

$$\text{Err}(\widehat{\Psi}) \xrightarrow{\text{prob.}} \Phi\left(-\frac{\gamma^2}{2\sqrt{\gamma^2 + 2\alpha}}\right)$$

What can help in high-dimensions?

What can help in high-dimensions?

Structure

What can help in high-dimensions?

Structure e.g., sparsity

What can help in high-dimensions?

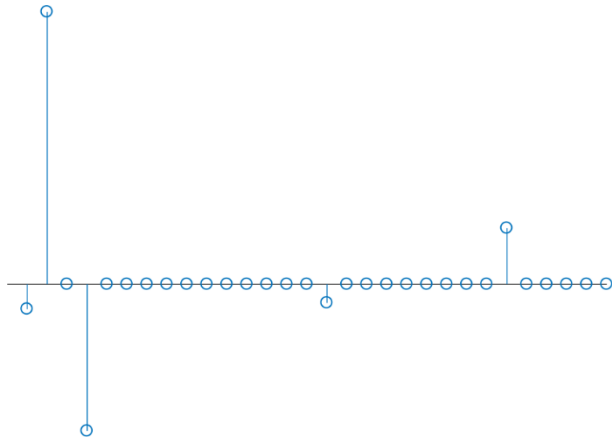
Structure e.g., sparsity

Suppose μ_0 and μ_1 are sparse: only have $s \ll d$ nonzero entries

What can help in high-dimensions?

Structure e.g., sparsity

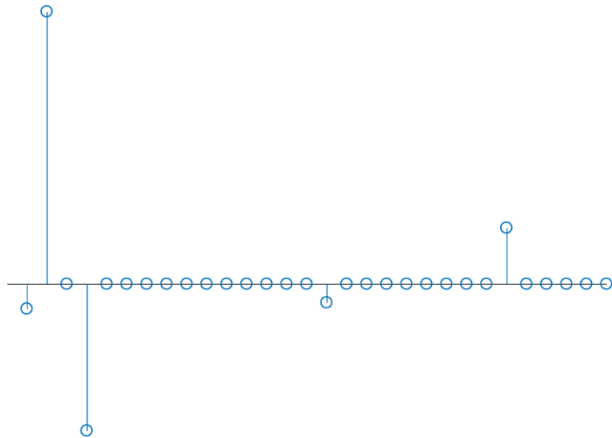
Suppose μ_0 and μ_1 are sparse: only have $s \ll d$ nonzero entries



What can help in high-dimensions?

Structure e.g., sparsity

Suppose μ_0 and μ_1 are sparse: only have $s \ll d$ nonzero entries

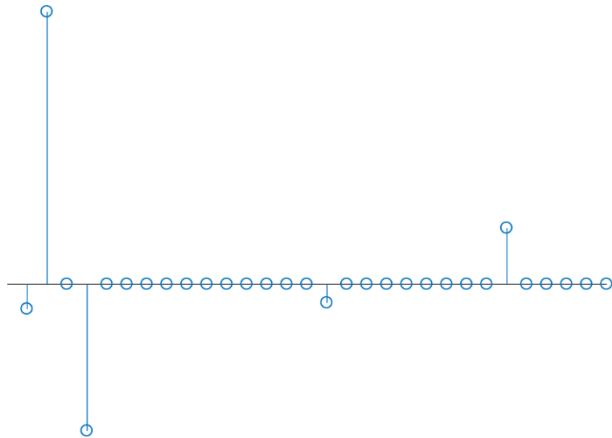


Procedure: *hard-threshold* entries of estimates

What can help in high-dimensions?

Structure e.g., sparsity

Suppose μ_0 and μ_1 are sparse: only have $s \ll d$ nonzero entries



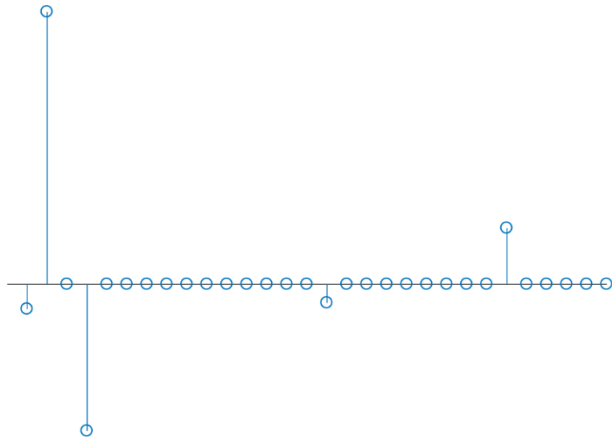
Procedure: *hard-threshold* entries of estimates

$$\hat{\mu}'_0 = \frac{1}{n_0} \sum_{i=1}^{n_0} x_i$$

What can help in high-dimensions?

Structure e.g., sparsity

Suppose μ_0 and μ_1 are sparse: only have $s \ll d$ nonzero entries



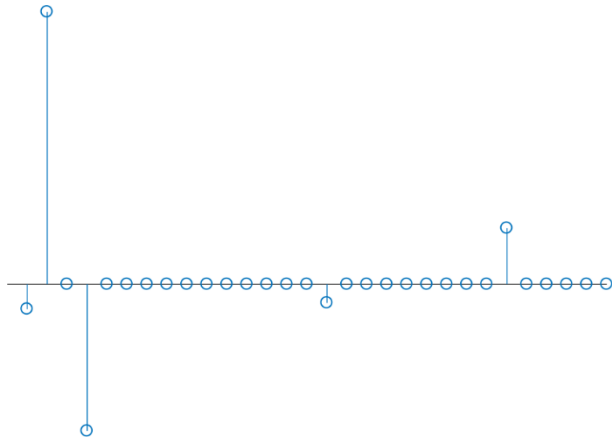
Procedure: *hard-threshold* entries of estimates

$$\hat{\mu}'_0 = \frac{1}{n_0} \sum_{i=1}^{n_0} x_i \quad \longrightarrow \quad (\hat{\mu}_0)_i = H_\lambda \left((\hat{\mu}'_0)_i \right)$$

What can help in high-dimensions?

Structure e.g., sparsity

Suppose μ_0 and μ_1 are sparse: only have $s \ll d$ nonzero entries



Procedure: *hard-threshold* entries of estimates

$$\hat{\mu}'_0 = \frac{1}{n_0} \sum_{i=1}^{n_0} x_i \quad \longrightarrow \quad (\hat{\mu}_0)_i = H_\lambda \left((\hat{\mu}'_0)_i \right)$$

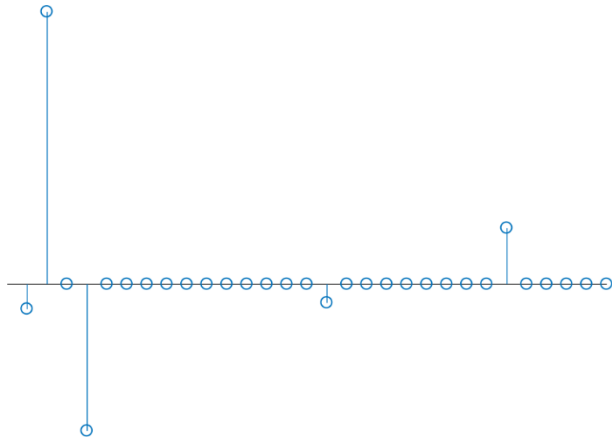
hard-thresholding operator

$$H_\lambda(x) := \begin{cases} x & , \text{ if } |x| > \lambda \\ 0 & , \text{ if } |x| \leq \lambda \end{cases}$$

What can help in high-dimensions?

Structure e.g., sparsity

Suppose μ_0 and μ_1 are sparse: only have $s \ll d$ nonzero entries



Procedure: *hard-threshold* entries of estimates

$$\hat{\mu}'_0 = \frac{1}{n_0} \sum_{i=1}^{n_0} x_i \quad \longrightarrow \quad (\hat{\mu}_0)_i = H_\lambda \left((\hat{\mu}'_0)_i \right)$$

(same for μ_1)

hard-thresholding operator

$$H_\lambda(x) := \begin{cases} x & , \text{ if } |x| > \lambda \\ 0 & , \text{ if } |x| \leq \lambda \end{cases}$$

Example

hard-thresholding operator

$$H_{\lambda}(x) := \begin{cases} x & , \text{ if } |x| > \lambda \\ 0 & , \text{ if } |x| \leq \lambda \end{cases}$$

Example

$$d = 400 \quad n = 800 \quad s = 5 \quad \lambda = \sqrt{\frac{2 \log d}{n}} = 0.1224$$

hard-thresholding operator

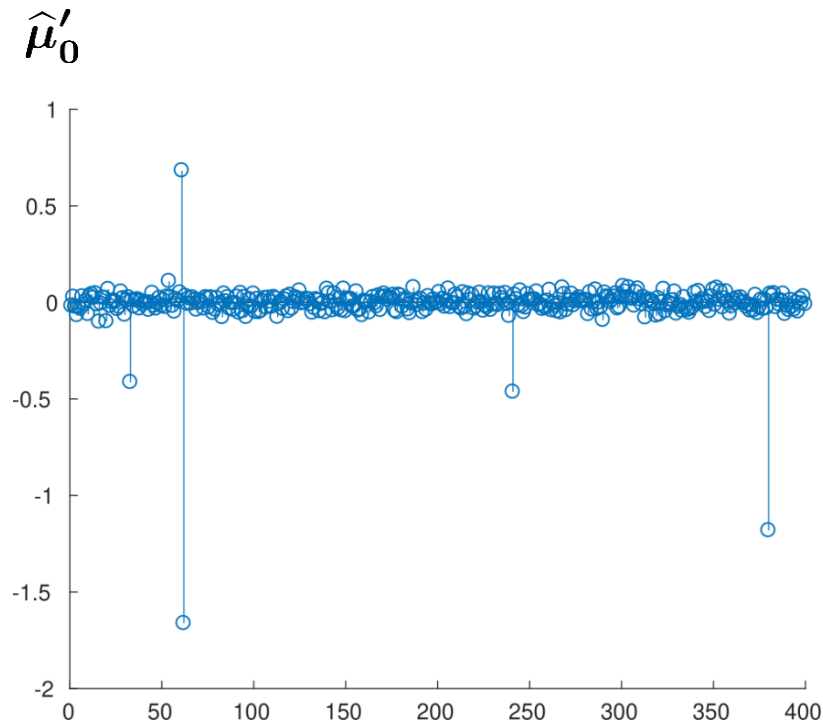
$$H_\lambda(x) := \begin{cases} x & , \text{ if } |x| > \lambda \\ 0 & , \text{ if } |x| \leq \lambda \end{cases}$$

Example

$$d = 400 \quad n = 800 \quad s = 5 \quad \lambda = \sqrt{\frac{2 \log d}{n}} = 0.1224$$

hard-thresholding operator

$$H_\lambda(x) := \begin{cases} x & , \text{ if } |x| > \lambda \\ 0 & , \text{ if } |x| \leq \lambda \end{cases}$$

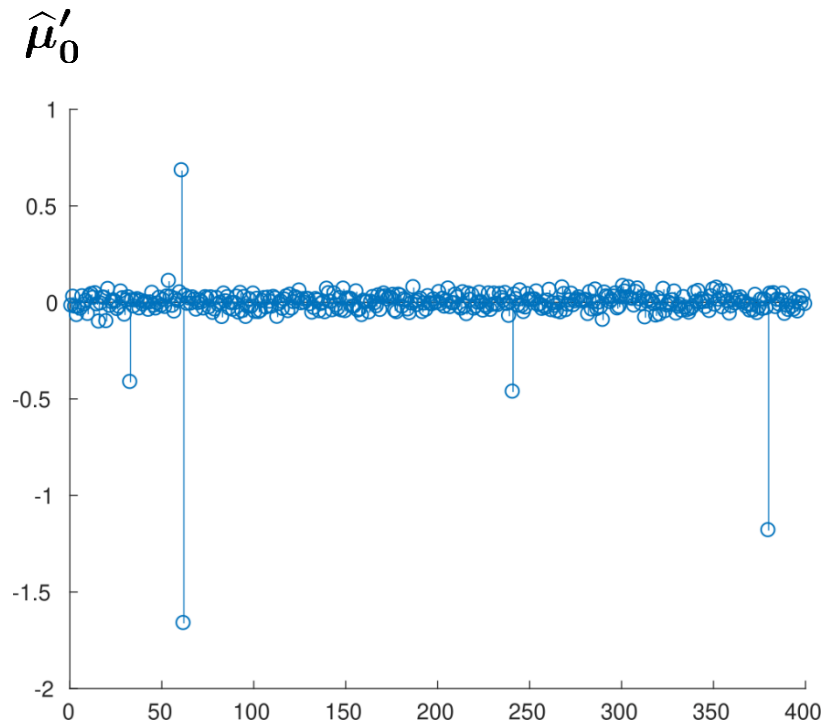


Example

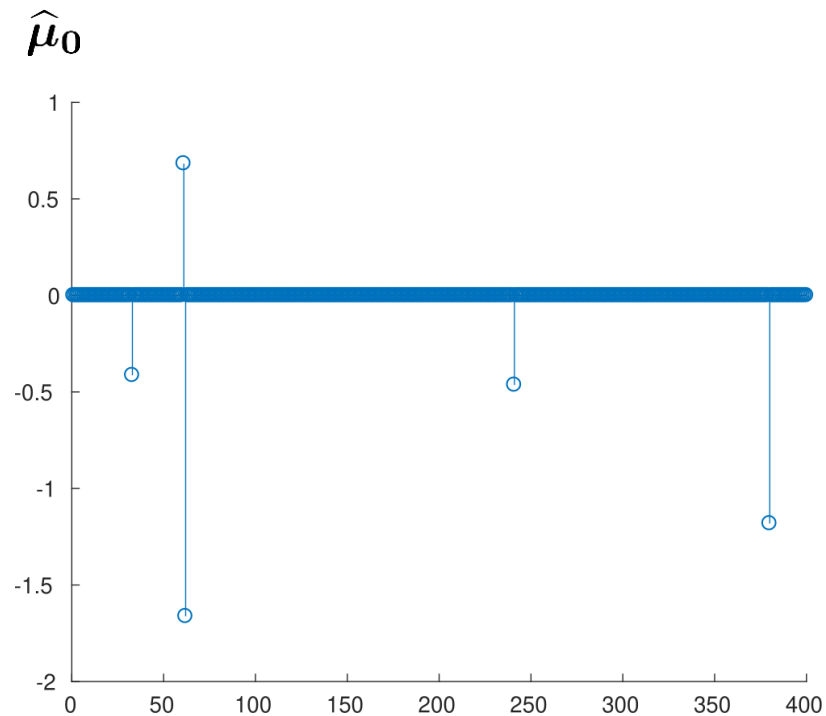
$$d = 400 \quad n = 800 \quad s = 5 \quad \lambda = \sqrt{\frac{2 \log d}{n}} = 0.1224$$

hard-thresholding operator

$$H_\lambda(x) := \begin{cases} x & , \text{ if } |x| > \lambda \\ 0 & , \text{ if } |x| \leq \lambda \end{cases}$$



$\xrightarrow{H_\lambda(\cdot)}$



Same Experiments

Same Experiments

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

$d = 400$ $n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\hat{\Psi}(x)$ over 5000 random trials

Same Experiments

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

$d = 400$ $n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\widehat{\Psi}(x)$ over 5000 random trials

$$s = 5 \quad \lambda = \sqrt{\frac{2 \log d}{n}}$$

Same Experiments

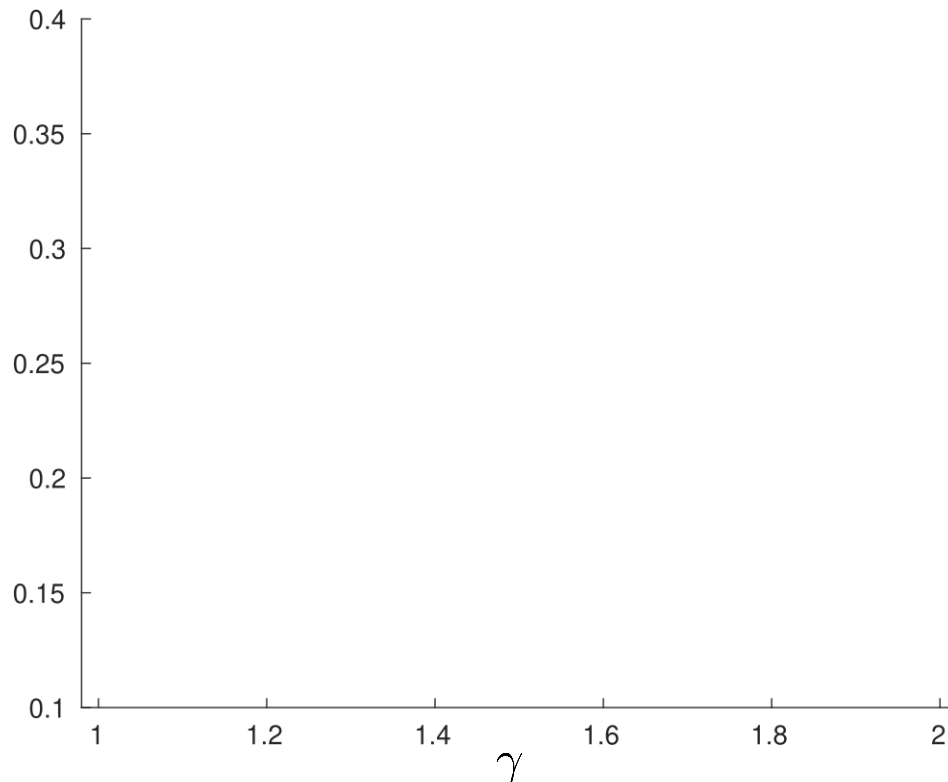
Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

$d = 400$ $n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\hat{\Psi}(x)$ over 5000 random trials

Probability of error



$$s = 5 \quad \lambda = \sqrt{\frac{2 \log d}{n}}$$

Same Experiments

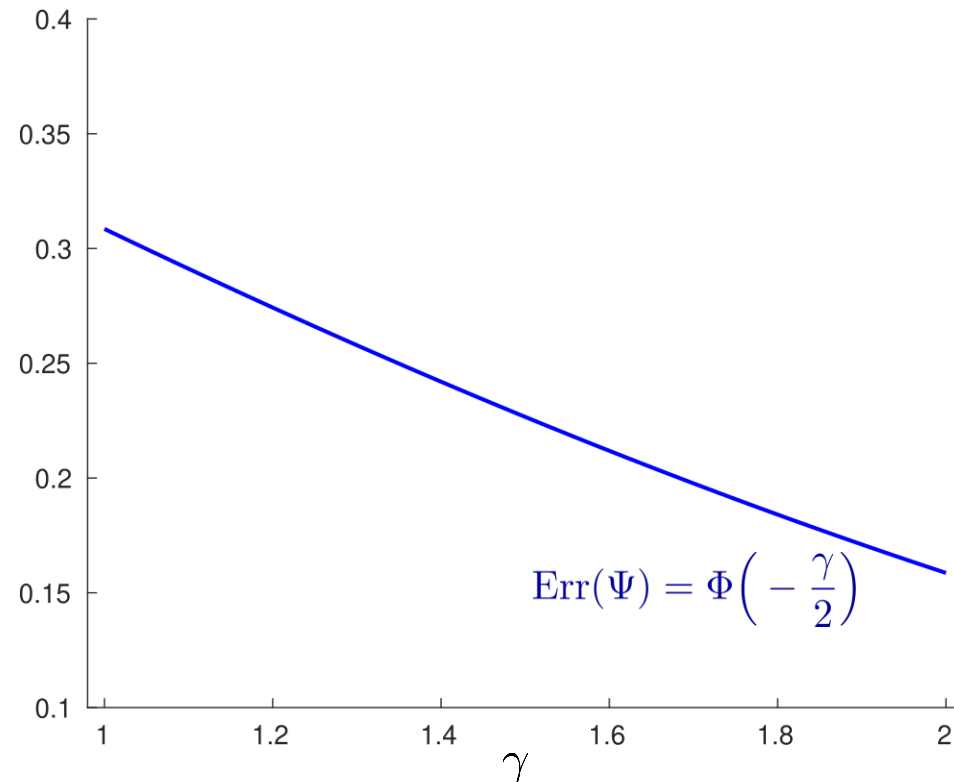
Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

$d = 400$ $n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\widehat{\Psi}(x)$ over 5000 random trials

Probability of error



$$s = 5 \quad \lambda = \sqrt{\frac{2 \log d}{n}}$$

Same Experiments

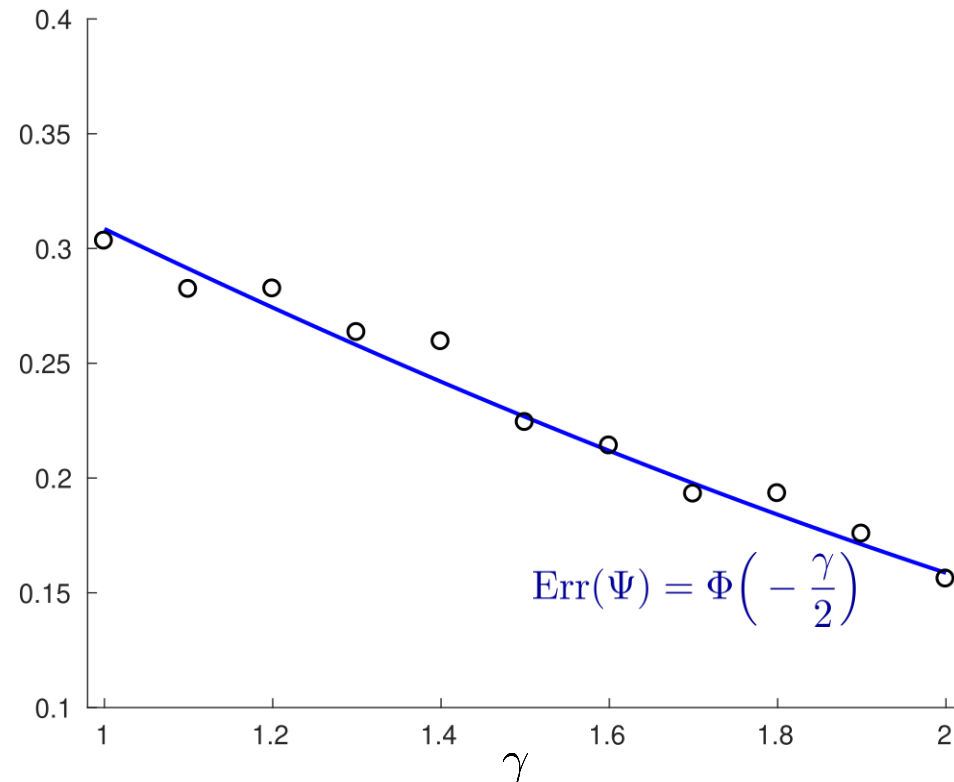
Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

$d = 400$ $n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\hat{\Psi}(x)$ over 5000 random trials

Probability of error



$$s = 5 \quad \lambda = \sqrt{\frac{2 \log d}{n}}$$

Same Experiments

Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

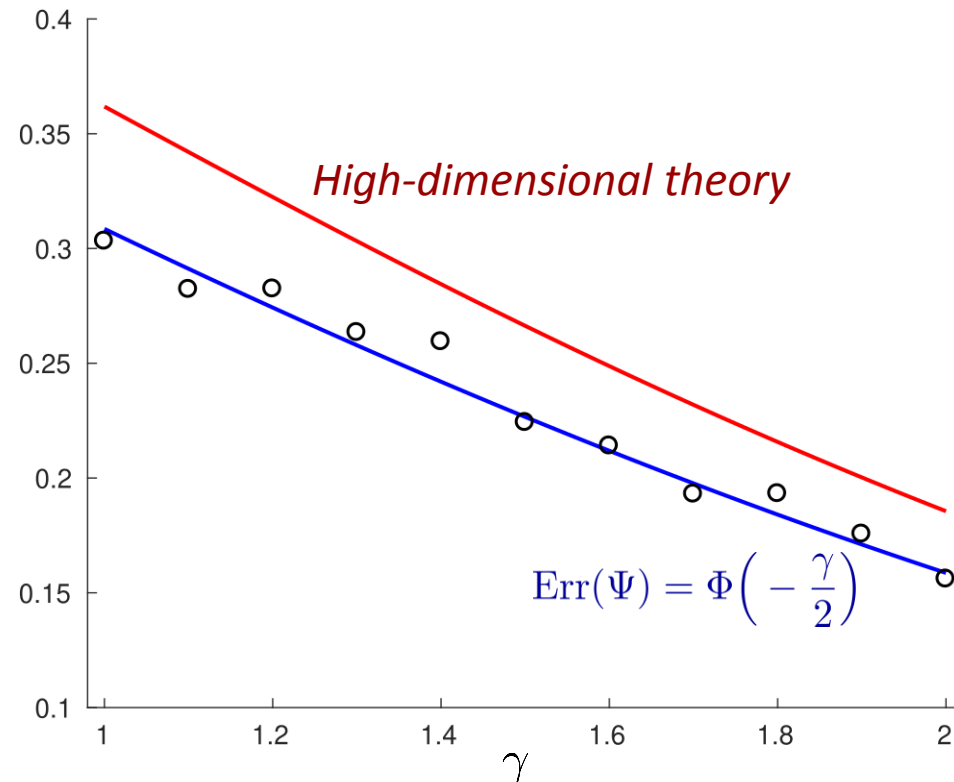
$d = 400$ $n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\hat{\Psi}(x)$ over 5000 random trials

$$s = 5 \quad \lambda = \sqrt{\frac{2 \log d}{n}}$$

Probability of error



Same Experiments

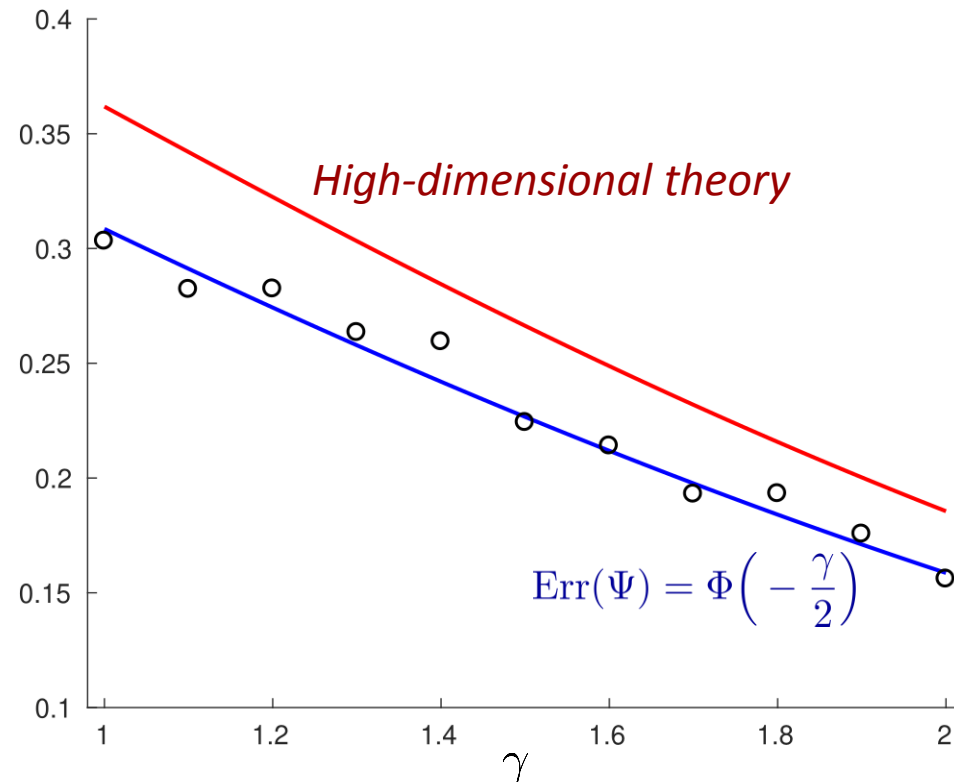
Assume $\mathbb{P}(\mathbb{P}_0) = \mathbb{P}(\mathbb{P}_1)$ $\Sigma = I_d$

$d = 400$ $n_0 = n_1 = 800$

Vary $\gamma = \|\mu_0 - \mu_1\|_2$ between 1 and 2

Test with $\hat{\Psi}(x)$ over 5000 random trials

Probability of error



$$s = 5 \quad \lambda = \sqrt{\frac{2 \log d}{n}}$$

Sparsity makes problem low-dimensional

Outline

- *Motivation: Hypothesis Testing in High-Dimensions*
- Introduction to LASSO and other sparsity problems
- Gaussian graphical model selection
- Matrix completion

A Crime Problem

A Crime Problem

City	Police funding/ resident (\$/ya)	% with 4 years high-school (+25)	% not in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes/million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

A Crime Problem

City	Police funding/ resident (\$/ya)	% with 4 years high-school (+25)	% not in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes/million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

Goal: Predict *# crimes / million* based on the other indicators

Linear Regression

City	Police fundings per resident	% with 4 years high-school (+25)	% in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes / million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

Linear Regression

Linear/affine model

City	Police fundings per resident	% with 4 years high-school (+25)	% in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes /million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

Linear Regression

Linear/affine model

$d = 5$ predictors

City	Police fundings per resident	% with 4 years high-school (+25)	% in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes / million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

← d →

Linear Regression

Linear/affine model

$d = 5$ predictors

$n = 50$ samples

City	Police fundings per resident	% with 4 years high-school (+25)	% in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes /million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

Linear Regression

Linear/affine model

$d = 5$ predictors

$n = 50$ samples

$$y_i \simeq x_0 + \sum_{j=1}^d a_{ij} x_j \quad i = 1, \dots, n$$

City	Police fundings per resident	% with 4 years high-school (+25)	% in high-school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes /million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

Linear Regression

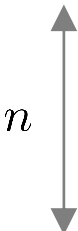
Linear/affine model

$d = 5$ predictors

$n = 50$ samples

$$y_i \simeq x_0 + \sum_{j=1}^d a_{ij} x_j \quad i = 1, \dots, n$$

*response variable
(crime rate)*



City	Police funding per resident	% with 4+ years high-school (+25)	% in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes / million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

Linear Regression

Linear/affine model

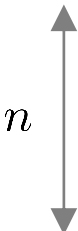
$d = 5$ predictors

$n = 50$ samples

$$y_i \simeq x_0 + \sum_{j=1}^d a_{ij} x_j \quad i = 1, \dots, n$$

offset

*response variable
(crime rate)*



City	Police fundings per resident	% with 4 years high-school (+25)	% in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes / million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

Linear Regression

Linear/affine model

$d = 5$ predictors

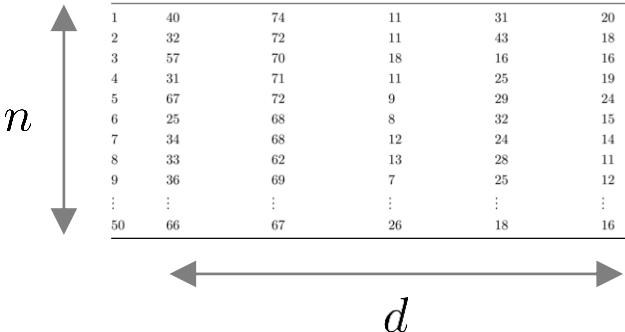
$n = 50$ samples

$$y_i \simeq x_0 + \sum_{j=1}^d a_{ij} x_j \quad i = 1, \dots, n$$

offset

*response variable
(crime rate)*

table entry



City	Police fundings per resident	% with 4+ years high-school (+25)	% in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes / million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

Linear Regression

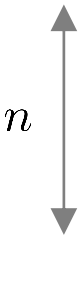
Linear/affine model

$d = 5$ predictors

$n = 50$ samples

$$y_i \simeq x_0 + \sum_{j=1}^d a_{ij} x_j \quad i = 1, \dots, n$$

y_i : response variable (crime rate)
 x_0 : offset
 a_{ij} : coefficient to be determined
 x_j : table entry



City	Police fundings per resident	% with 4+ years high-school (+25)	% in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes / million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

Linear Regression

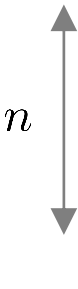
Linear/affine model

$d = 5$ predictors

$n = 50$ samples

$$y_i \simeq x_0 + \sum_{j=1}^d a_{ij} x_j \quad i = 1, \dots, n \quad \Longleftrightarrow$$

y_i : response variable (crime rate)
 x_0 : offset
 a_{ij} : coefficient to be determined
 x_j : table entry



City	Police fundings per resident	% with 4 years high-school (+25)	% in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes / million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

d

$$y \simeq x_0 1_n + Ax$$

Linear Regression

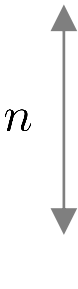
Linear/affine model

$d = 5$ predictors

$n = 50$ samples

$$y_i \simeq x_0 + \sum_{j=1}^d a_{ij} x_j \quad i = 1, \dots, n \quad \Longleftrightarrow$$

y_i : response variable (crime rate)
 x_0 : offset
 a_{ij} : coefficient to be determined
 x_j : table entry



City	Police fundings per resident	% with 4 years high-school (+25)	% in high-school (16-19)	% in college (18-24)	% with 4 years college (25+)	Crimes / million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

$$y \simeq x_0 1_n + Ax$$

$$= \begin{bmatrix} 1_n & A \end{bmatrix} \begin{bmatrix} x_0 \\ x \end{bmatrix}$$

Linear Regression

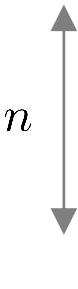
Linear/affine model

$d = 5$ predictors

$n = 50$ samples

$$y_i \simeq x_0 + \sum_{j=1}^d a_{ij} x_j \quad i = 1, \dots, n \quad \Longleftrightarrow$$

y_i : response variable (crime rate)
 x_0 : offset
 a_{ij} : coefficient to be determined
 x_j : table entry



City	Police fundings per resident	% with 4 years high-school (+25)	% in high-school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes / million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

$$\begin{aligned}
 y &\simeq x_0 1_n + Ax \\
 &= \underbrace{[1_n \quad A]}_{\overline{A}} \underbrace{\begin{bmatrix} x_0 \\ x \end{bmatrix}}_{\overline{x}}
 \end{aligned}$$

Linear Regression

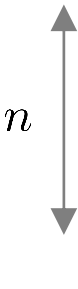
Linear/affine model

$d = 5$ predictors

$n = 50$ samples

$$y_i \simeq x_0 + \sum_{j=1}^d a_{ij} x_j \quad i = 1, \dots, n \quad \Longleftrightarrow$$

y_i : response variable (crime rate)
 x_0 : offset
 a_{ij} : coefficient to be determined
 x_j : table entry



City	Police fundings per resident	% with 4 years high-school (+25)	% in high-school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes / million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

$$y \simeq x_0 1_n + Ax$$

$$= \underbrace{[1_n \quad A]}_{\bar{A}} \underbrace{\begin{bmatrix} x_0 \\ x \end{bmatrix}}_{\bar{x}}$$

Find coefficients:

Linear Regression

Linear/affine model

$d = 5$ predictors

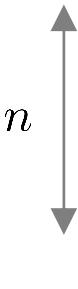
$n = 50$ samples

$$y_i \simeq x_0 + \sum_{j=1}^d a_{ij} x_j \quad i = 1, \dots, n \quad \Longleftrightarrow$$

y_i : response variable (crime rate)
 x_0 : offset
 a_{ij} : coefficient to be determined
 x_j : table entry

Find coefficients: *least-squares*

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \|y - \bar{A}\bar{x}\|_2^2$$



City	Police funding per resident	% with 4 years high-school (+25)	% in high-school (16-19)	% in college (18-24)	% with 4 years college (25+)	Crimes / million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

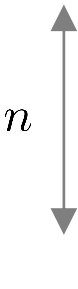
$$\begin{aligned}
 y &\simeq x_0 1_n + Ax \\
 &= \underbrace{[1_n \quad A]}_{\bar{A}} \underbrace{\begin{bmatrix} x_0 \\ x \end{bmatrix}}_{\bar{x}}
 \end{aligned}$$

Linear Regression

Linear/affine model

$d = 5$ predictors

$n = 50$ samples



City	Police funding per resident	% with 4 years high-school (+25)	% in high- school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Crimes /million
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
⋮	⋮	⋮	⋮	⋮	⋮	⋮
50	66	67	26	18	16	940

$$y_i \simeq x_0 + \sum_{j=1}^d a_{ij} x_j \quad i = 1, \dots, n \quad \Longleftrightarrow$$

y_i
 response variable
 (crime rate)

x_0
 offset

a_{ij}
 coefficient to be determined

x_j
 table entry

$$\begin{aligned}
 y &\simeq x_0 \mathbf{1}_n + A x \\
 &= \underbrace{[\mathbf{1}_n \quad A]}_{\overline{A}} \underbrace{\begin{bmatrix} x_0 \\ x \end{bmatrix}}_{\overline{x}}
 \end{aligned}$$

Find coefficients: *least-squares*

$$\underset{\overline{x}}{\text{minimize}} \quad \frac{1}{2} \left\| y - \overline{A} \overline{x} \right\|_2^2$$

\Longleftrightarrow

$$\overline{x}_{\text{LS}}^* = \left(\overline{A}^\top \overline{A} \right)^{-1} \overline{A} y$$

Linear Regression

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \left\| y - \overline{A} \bar{x} \right\|_2^2$$

$$\bar{x}_{\text{LS}}^* = \left(\overline{A}^\top \overline{A} \right)^{-1} \overline{A} y$$

Linear Regression

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \|y - \bar{A}\bar{x}\|_2^2$$

$$\bar{x}_{\text{LS}}^* = \left(\bar{A}^\top \bar{A}\right)^{-1} \bar{A} y$$

$$= \begin{bmatrix} 489.6486 \\ 10.9807 \\ -6.0885 \\ 5.4803 \\ 0.3770 \\ 5.5005 \end{bmatrix}$$

Linear Regression

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \|y - \bar{A}\bar{x}\|_2^2$$

$$\bar{x}_{\text{LS}}^* = \left(\bar{A}^\top \bar{A}\right)^{-1} \bar{A} y$$

$$= \begin{bmatrix} 489.6486 \\ 10.9807 \\ -6.0885 \\ 5.4803 \\ 0.3770 \\ 5.5005 \end{bmatrix} \text{ offset}$$

Linear Regression

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \|y - \bar{A}\bar{x}\|_2^2$$

$$\bar{x}_{\text{LS}}^* = \left(\bar{A}^\top \bar{A}\right)^{-1} \bar{A} y$$

$$= \begin{bmatrix} 489.6486 \\ 10.9807 \\ -6.0885 \\ 5.4803 \\ 0.3770 \\ 5.5005 \end{bmatrix} \begin{array}{l} \textit{offset} \\ \textit{police funding / resident (\$/year)} \end{array}$$

Linear Regression

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \|y - \bar{A}\bar{x}\|_2^2$$

$$\bar{x}_{\text{LS}}^* = (\bar{A}^\top \bar{A})^{-1} \bar{A} y$$

$$= \begin{bmatrix} 489.6486 \\ 10.9807 \\ -6.0885 \\ 5.4803 \\ 0.3770 \\ 5.5005 \end{bmatrix} \begin{array}{l} \textit{offset} \\ \textit{police funding / resident (\$/year)} \\ \textit{\% of 25+ year-olds with 4+ years of high-school} \\ \\ \\ \end{array}$$

Linear Regression

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \|y - \bar{A}\bar{x}\|_2^2$$

$$\bar{x}_{\text{LS}}^* = (\bar{A}^\top \bar{A})^{-1} \bar{A} y$$

$$= \begin{bmatrix} 489.6486 \\ 10.9807 \\ -6.0885 \\ 5.4803 \\ 0.3770 \\ 5.5005 \end{bmatrix} \begin{array}{l} \textit{offset} \\ \textit{police funding / resident (\$/year)} \\ \textit{\% of 25+ year-olds with 4+ years of high-school} \\ \textit{\% of 16-19 year-olds not in high-school} \\ \textit{\% of 18-24 year-olds in college} \\ \textit{\% of 25+ year-olds with 4+ years of college} \end{array}$$

Linear Regression

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \|y - \bar{A}\bar{x}\|_2^2$$

$$\bar{x}_{\text{LS}}^* = (\bar{A}^\top \bar{A})^{-1} \bar{A} y$$

$$= \begin{bmatrix} 489.6486 \\ 10.9807 \\ -6.0885 \\ 5.4803 \\ 0.3770 \\ 5.5005 \end{bmatrix} \begin{array}{l} \textit{offset} \\ \textit{police funding / resident (\$/year)} \\ \textit{\% of 25+ year-olds with 4+ years of high-school} \\ \textit{\% of 16-19 year-olds not in high-school} \\ \textit{\% of 18-24 year-olds in college} \\ \textit{\% of 25+ year-olds with 4+ years of college} \end{array}$$

Problems with least-squares

Linear Regression

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \|y - \bar{A}\bar{x}\|_2^2$$

$$\bar{x}_{\text{LS}}^* = (\bar{A}^\top \bar{A})^{-1} \bar{A} y$$

$$= \begin{bmatrix} 489.6486 \\ 10.9807 \\ -6.0885 \\ 5.4803 \\ 0.3770 \\ 5.5005 \end{bmatrix} \begin{array}{l} \text{offset} \\ \text{police funding / resident (\$/year)} \\ \text{\% of 25+ year-olds with 4+ years of high-school} \\ \text{\% of 16-19 year-olds not in high-school} \\ \text{\% of 18-24 year-olds in college} \\ \text{\% of 25+ year-olds with 4+ years of college} \end{array}$$

Problems with least-squares

little interpretability

Linear Regression

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \|y - \bar{A}\bar{x}\|_2^2$$

$$\bar{x}_{\text{LS}}^* = (\bar{A}^\top \bar{A})^{-1} \bar{A} y$$

$$= \begin{bmatrix} 489.6486 \\ 10.9807 \\ -6.0885 \\ 5.4803 \\ 0.3770 \\ 5.5005 \end{bmatrix} \begin{array}{l} \text{offset} \\ \text{police funding / resident (\$/year)} \\ \text{\% of 25+ year-olds with 4+ years of high-school} \\ \text{\% of 16-19 year-olds not in high-school} \\ \text{\% of 18-24 year-olds in college} \\ \text{\% of 25+ year-olds with 4+ years of college} \end{array}$$

Problems with least-squares

little interpretability

all coefficients contribute to prediction

Linear Regression

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \|y - \bar{A}\bar{x}\|_2^2$$

$$\bar{x}_{\text{LS}}^* = (\bar{A}^\top \bar{A})^{-1} \bar{A} y$$

$$= \begin{bmatrix} 489.6486 \\ 10.9807 \\ -6.0885 \\ 5.4803 \\ 0.3770 \\ 5.5005 \end{bmatrix} \begin{array}{l} \text{offset} \\ \text{police funding / resident (\$/year)} \\ \text{\% of 25+ year-olds with 4+ years of high-school} \\ \text{\% of 16-19 year-olds not in high-school} \\ \text{\% of 18-24 year-olds in college} \\ \text{\% of 25+ year-olds with 4+ years of college} \end{array}$$

Problems with least-squares

little interpretability

all coefficients contribute to prediction

small bias, large variance

Linear Regression

$$\underset{\bar{x}}{\text{minimize}} \quad \frac{1}{2} \|y - \bar{A}\bar{x}\|_2^2$$

$$\bar{x}_{\text{LS}}^* = (\bar{A}^\top \bar{A})^{-1} \bar{A} y$$

$$= \begin{bmatrix} 489.6486 \\ 10.9807 \\ -6.0885 \\ 5.4803 \\ 0.3770 \\ 5.5005 \end{bmatrix} \begin{array}{l} \text{offset} \\ \text{police funding / resident (\$/year)} \\ \text{\% of 25+ year-olds with 4+ years of high-school} \\ \text{\% of 16-19 year-olds not in high-school} \\ \text{\% of 18-24 year-olds in college} \\ \text{\% of 25+ year-olds with 4+ years of college} \end{array}$$

Problems with least-squares

little interpretability

all coefficients contribute to prediction

small bias, large variance

zeroing coefficients can improve mean-squared error

LASSO

LASSO

least absolute selection and shrinkage operator

LASSO

least absolute selection and shrinkage operator

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

LASSO

least absolute selection and shrinkage operator

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

| regularization parameter ≥ 0

LASSO

least absolute selection and shrinkage operator

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

— *L1-norm*

| regularization parameter ≥ 0

LASSO

least absolute selection and shrinkage operator

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

———— **L1-norm** $\|x\|_1 = |x_1| + |x_2| + \cdots + |x_d|$

| regularization parameter ≥ 0

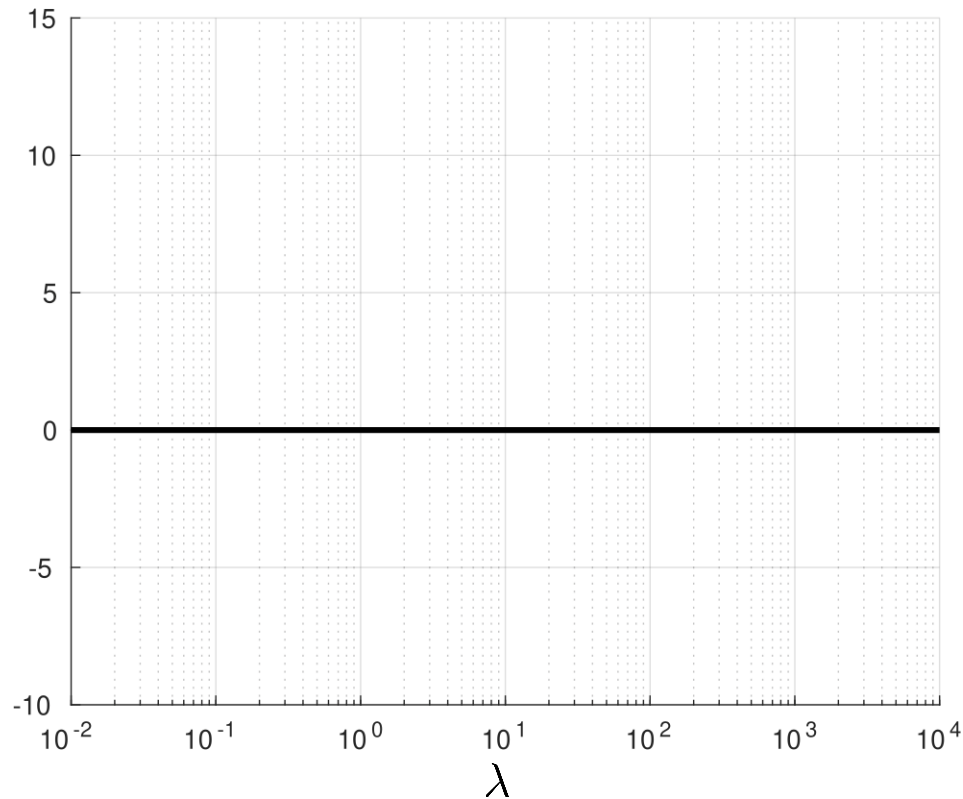
LASSO

least absolute selection and shrinkage operator

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1 \quad \text{--- } L1\text{-norm} \quad \|x\|_1 = |x_1| + |x_2| + \cdots + |x_d|$$

| regularization parameter ≥ 0

Coefficient value



LASSO

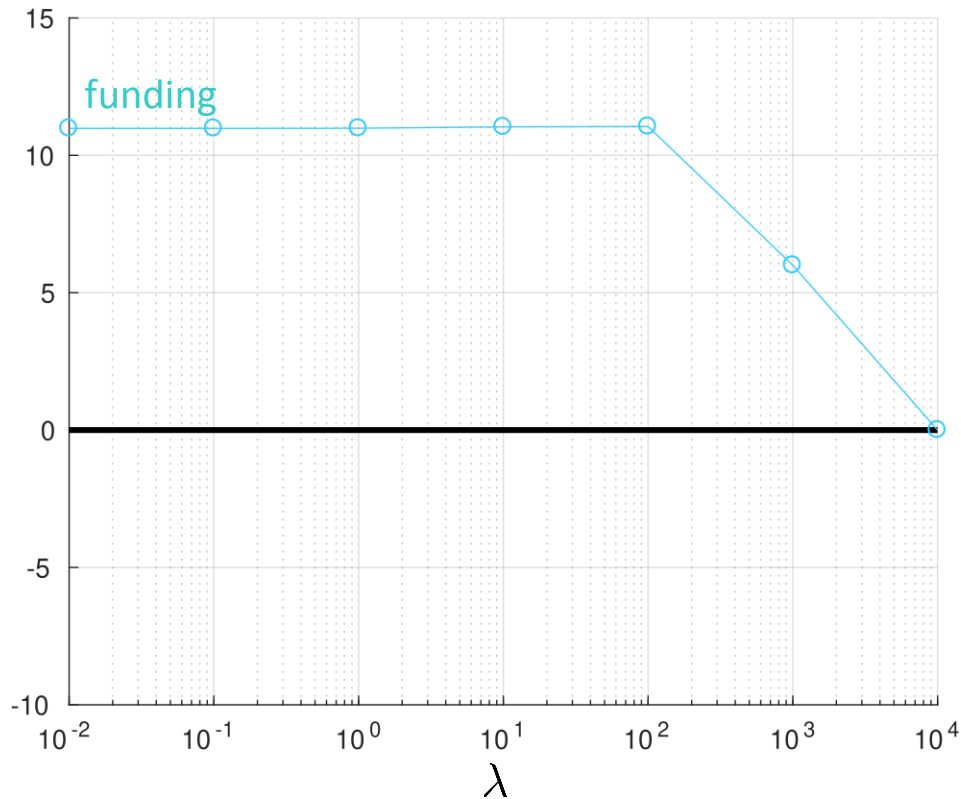
least absolute selection and shrinkage operator

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

———— **L1-norm** $\|x\|_1 = |x_1| + |x_2| + \cdots + |x_d|$

regularization parameter ≥ 0

Coefficient value



LASSO

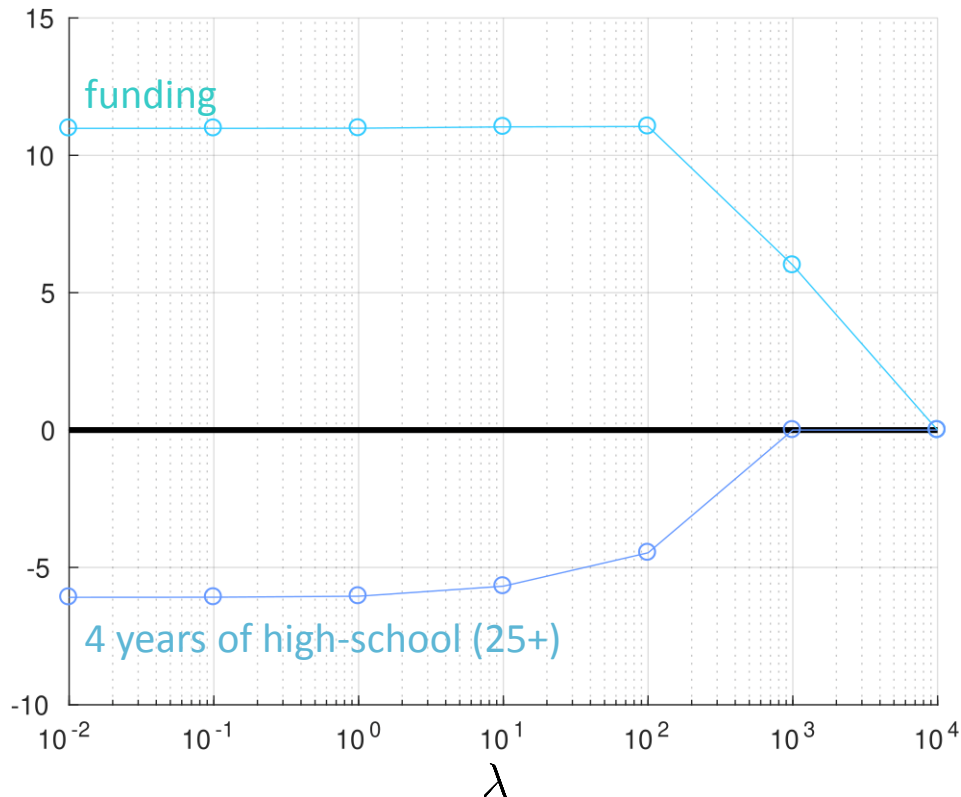
least absolute selection and shrinkage operator

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

———— **L1-norm** $\|x\|_1 = |x_1| + |x_2| + \cdots + |x_d|$

regularization parameter ≥ 0

Coefficient value



LASSO

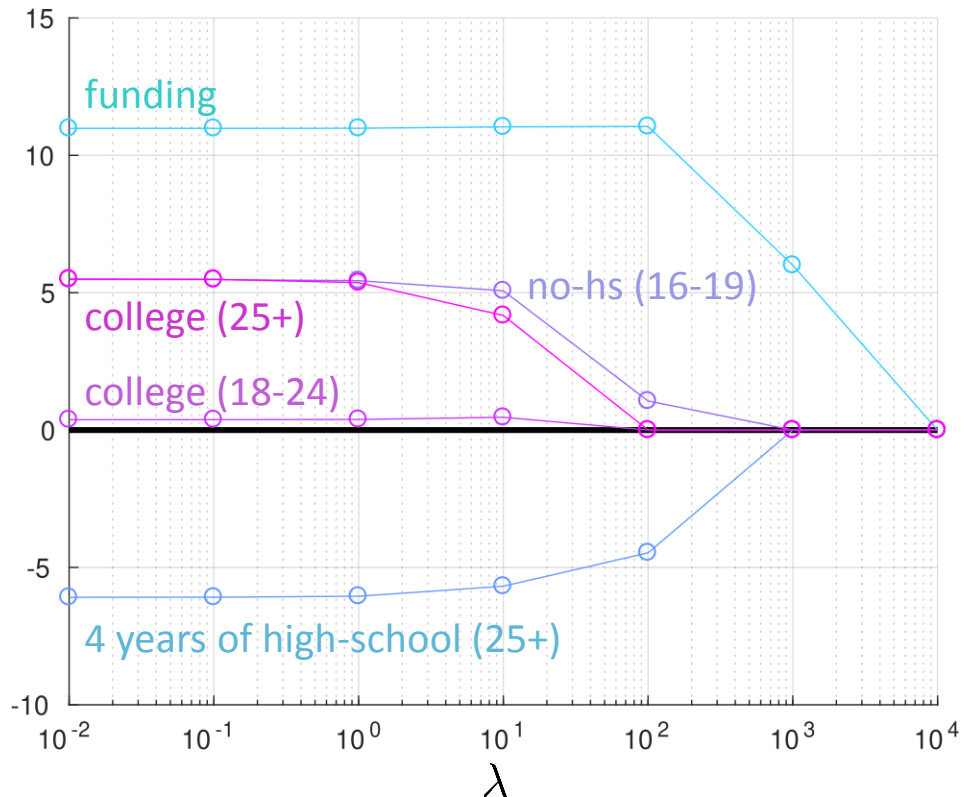
least absolute selection and shrinkage operator

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

L1-norm $\|x\|_1 = |x_1| + |x_2| + \dots + |x_d|$

regularization parameter ≥ 0

Coefficient value



LASSO

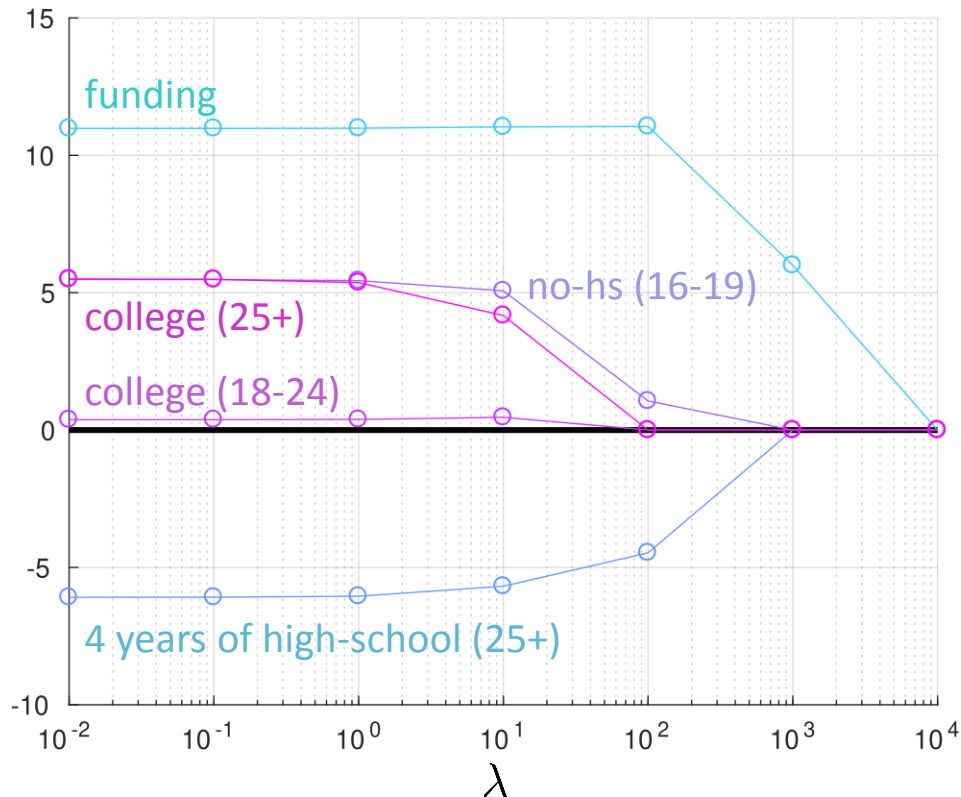
least absolute selection and shrinkage operator

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

L1-norm $\|x\|_1 = |x_1| + |x_2| + \dots + |x_d|$

regularization parameter ≥ 0

Coefficient value



In reality, we solved ...

LASSO

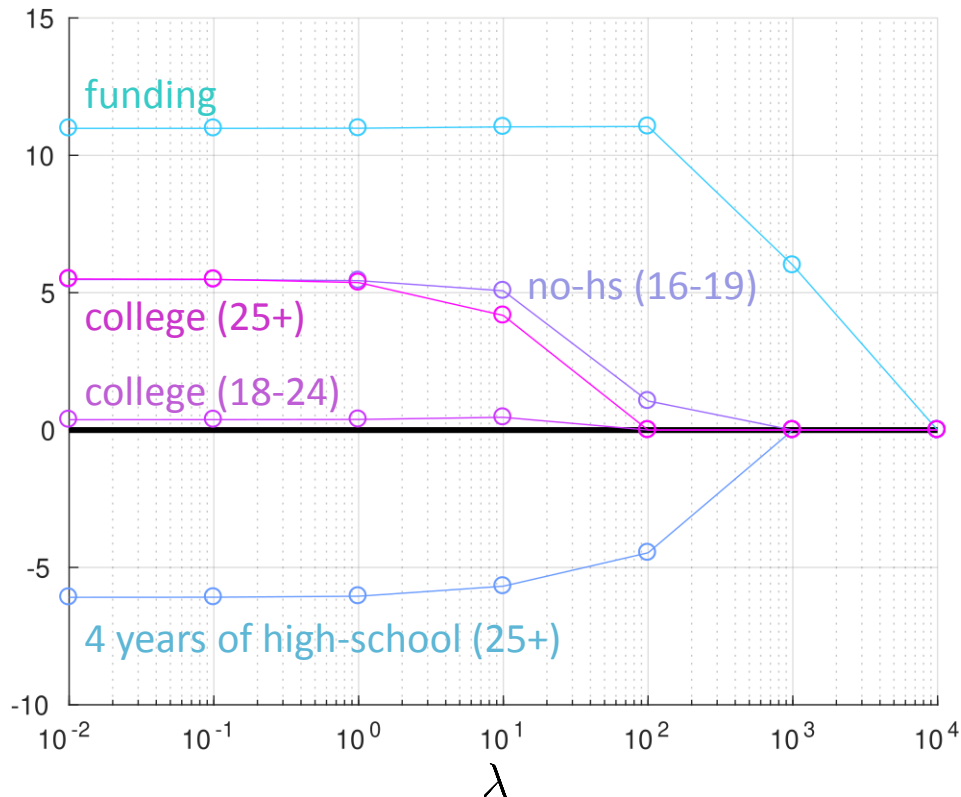
least absolute selection and shrinkage operator

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

L1-norm $\|x\|_1 = |x_1| + |x_2| + \dots + |x_d|$

regularization parameter ≥ 0

Coefficient value



In reality, we solved ...

$$\min_{x_0, x} \quad \frac{1}{2n} \|y - x_0 1_n - Ax\|_2^2 + \lambda \|x\|_1$$

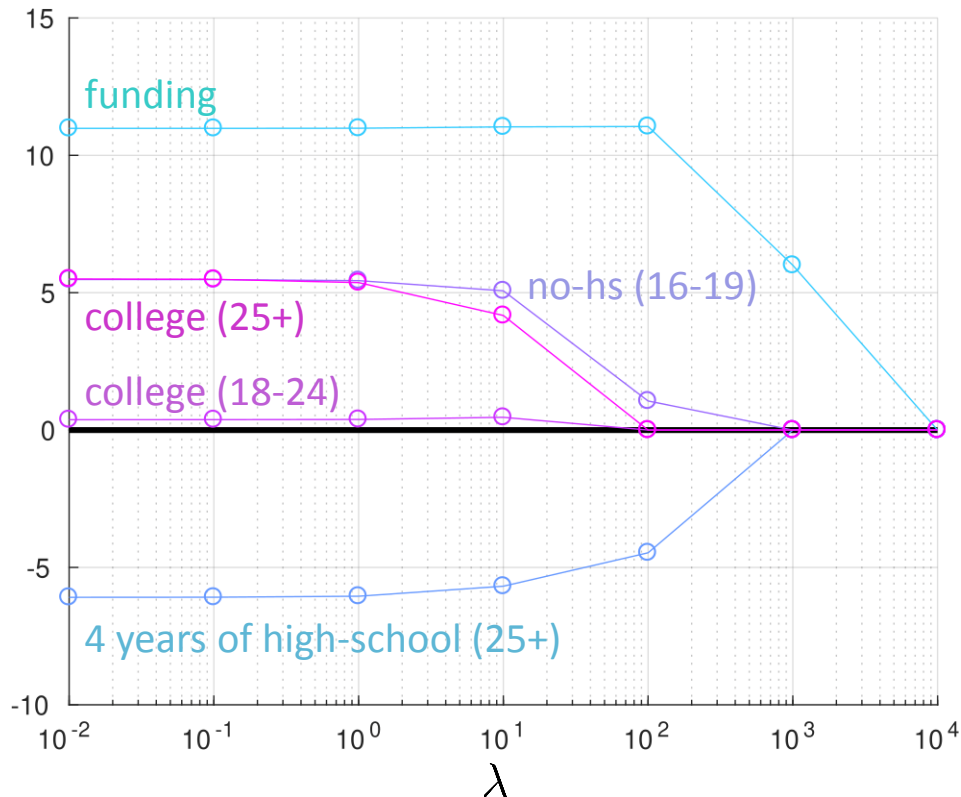
LASSO

least absolute selection and shrinkage operator

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1 \quad \text{--- } \textbf{L1-norm} \quad \|x\|_1 = |x_1| + |x_2| + \dots + |x_d|$$

| regularization parameter ≥ 0

Coefficient value



In reality, we solved ...

$$\min_{x_0, x} \quad \frac{1}{2n} \|y - x_0 1_n - Ax\|_2^2 + \lambda \|x\|_1$$

| necessary because $\frac{1}{n} \sum_{i=1}^n y_i \neq 0$

L1-Norm Induces Sparsity

L1-Norm Induces Sparsity

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

LASSO (aka Basis Pursuit Denoising)

L1-Norm Induces Sparsity

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

LASSO (aka Basis Pursuit Denoising)

\Updownarrow for some τ depending on λ

L1-Norm Induces Sparsity

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

LASSO (aka Basis Pursuit Denoising)

\Updownarrow for some τ depending on λ

$$\begin{array}{ll} \underset{x}{\text{minimize}} & \|y - Ax\|_2^2 \\ \text{subject to} & \|x\|_1 \leq \tau \end{array}$$

Constrained LASSO

L1-Norm Induces Sparsity

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

LASSO (aka Basis Pursuit Denoising)

\Updownarrow for some τ depending on λ

$$\begin{array}{ll} \underset{x}{\text{minimize}} & \|y - Ax\|_2^2 \\ \text{subject to} & \|x\|_1 \leq \tau \end{array}$$

Constrained LASSO

\Updownarrow for some σ depending on τ

L1-Norm Induces Sparsity

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

LASSO (aka Basis Pursuit Denoising)

\Updownarrow for some τ depending on λ

$$\begin{array}{ll} \underset{x}{\text{minimize}} & \|y - Ax\|_2^2 \\ \text{subject to} & \|x\|_1 \leq \tau \end{array}$$

Constrained LASSO

\Updownarrow for some σ depending on τ

$$\begin{array}{ll} \underset{x}{\text{minimize}} & \|x\|_1 \\ \text{subject to} & \|y - Ax\|_2 \leq \sigma \end{array}$$

Relaxed Basis Pursuit

L1-Norm Induces Sparsity

$$\underset{x}{\text{minimize}} \quad \frac{1}{2n} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

LASSO (aka Basis Pursuit Denoising)

\Updownarrow for some τ depending on λ

$$\begin{array}{ll} \underset{x}{\text{minimize}} & \|y - Ax\|_2^2 \\ \text{subject to} & \|x\|_1 \leq \tau \end{array}$$

Constrained LASSO

\Updownarrow for some σ depending on τ

$$\begin{array}{ll} \underset{x}{\text{minimize}} & \|x\|_1 \\ \text{subject to} & \|y - Ax\|_2 \leq \sigma \end{array}$$

Relaxed Basis Pursuit

Basis Pursuit when $\sigma = 0$

L1-Norm Induces Sparsity

$$\begin{aligned} \hat{x} \in \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & \|y - Ax\|_2 \leq \sigma \end{aligned}$$

L1-Norm Induces Sparsity

$$\begin{aligned} \hat{x} \in \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & \|y - Ax\|_2 \leq \sigma \end{aligned}$$

Assume $\sigma = 0$:

L1-Norm Induces Sparsity

$$\begin{aligned} \hat{x} \in \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & \|y - Ax\|_2 \leq \sigma \end{aligned}$$

Assume $\sigma = 0$:

$y = Ax$ has solutions

L1-Norm Induces Sparsity

$$\begin{aligned} \hat{x} \in \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & \|y - Ax\|_2 \leq \sigma \end{aligned}$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$

L1-Norm Induces Sparsity

$$\begin{aligned} \hat{x} \in \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & \|y - Ax\|_2 \leq \sigma \end{aligned}$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$
$$y = A\tilde{x}$$

L1-Norm Induces Sparsity

$$\begin{aligned} \hat{x} \in \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & \|y - Ax\|_2 \leq \sigma \end{aligned}$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$
$$y = A\tilde{x} \Big| \quad \Big| \{d : Ad = 0\}$$

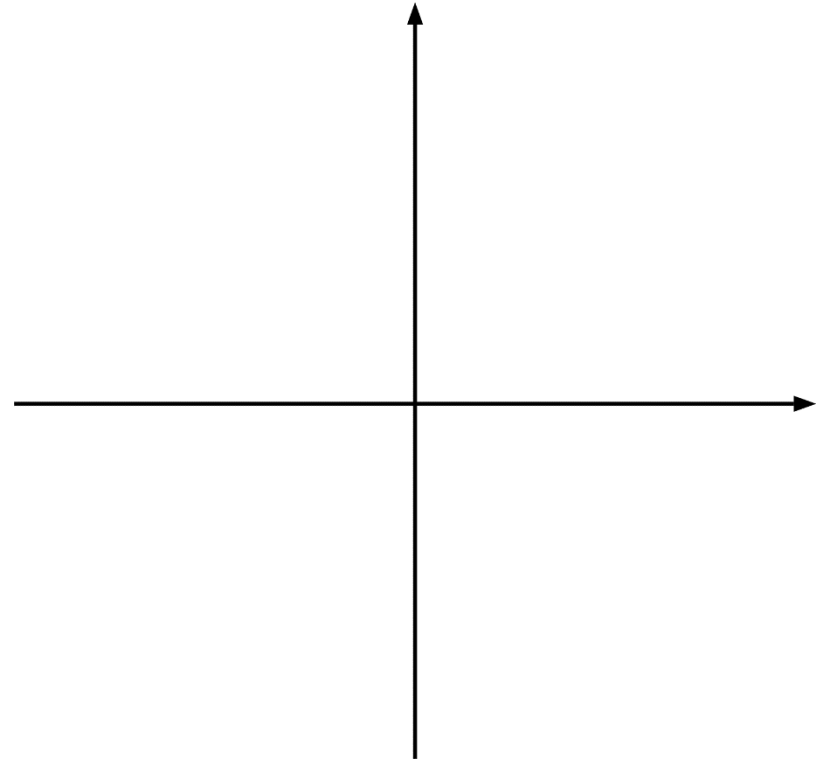
L1-Norm Induces Sparsity

$$\begin{aligned} \hat{x} \in \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & \|y - Ax\|_2 \leq \sigma \end{aligned}$$

Assume $\sigma = 0$:

$y = Ax$ has solutions $\tilde{x} + \text{null}(A)$

$$y = A\tilde{x} \Big| \quad \Big| \{d : Ad = 0\}$$



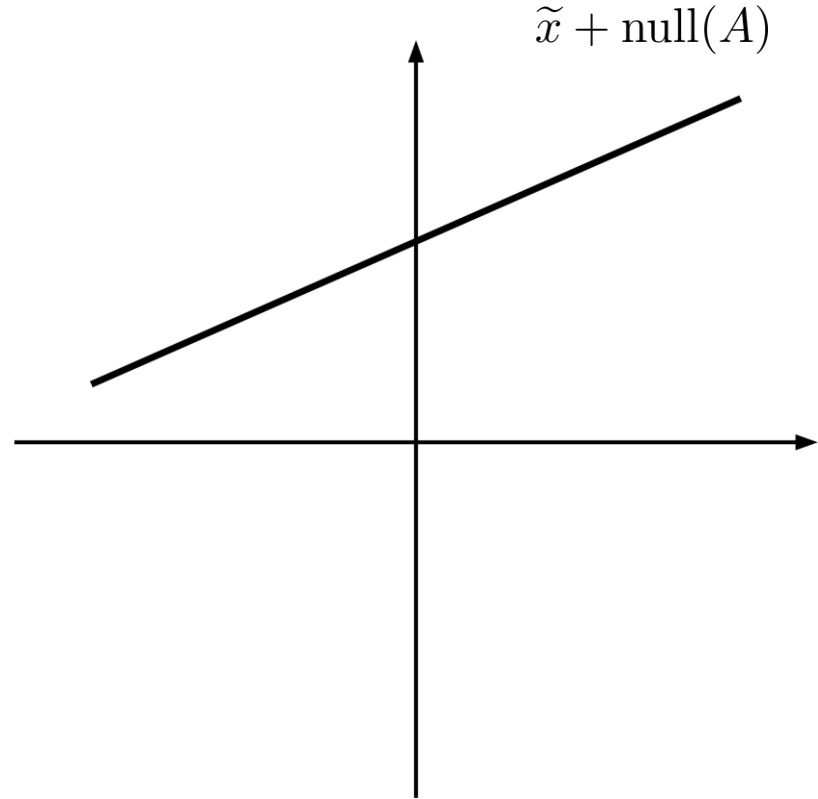
L1-Norm Induces Sparsity

$$\hat{x} \in \arg \min_x \|x\|_1$$
$$\text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$y = Ax$ has solutions $\tilde{x} + \text{null}(A)$

$$y = A\tilde{x} \Big| \quad \Big| \{d : Ad = 0\}$$



L1-Norm Induces Sparsity

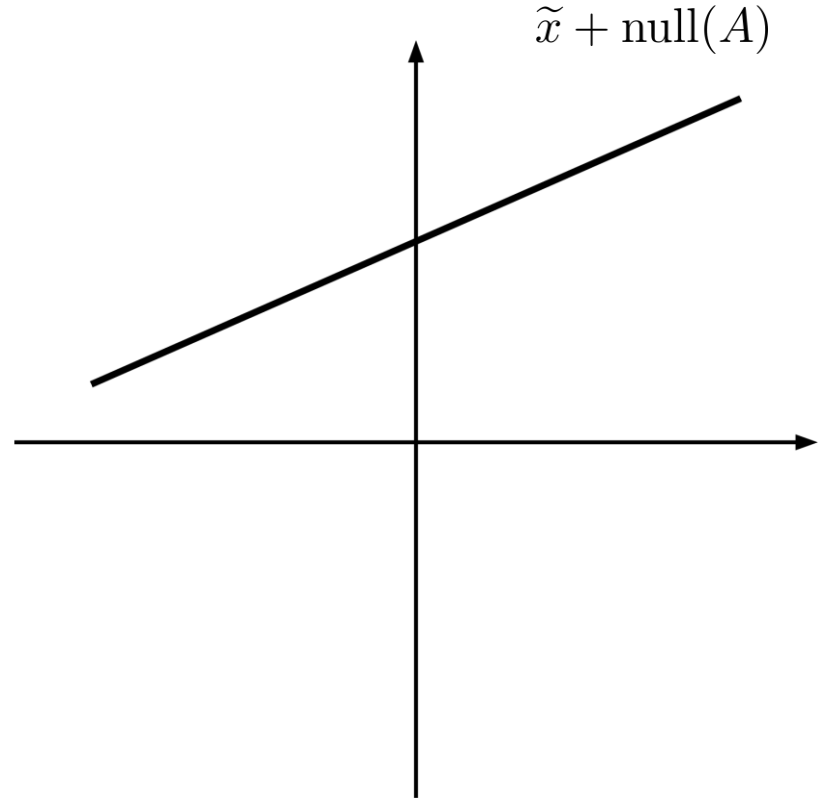
$$\hat{x} \in \arg \min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$y = Ax$ has solutions $\tilde{x} + \text{null}(A)$

$$y = A\tilde{x} \Big| \quad \Big| \{d : Ad = 0\}$$

Assume $\sigma > 0$:



L1-Norm Induces Sparsity

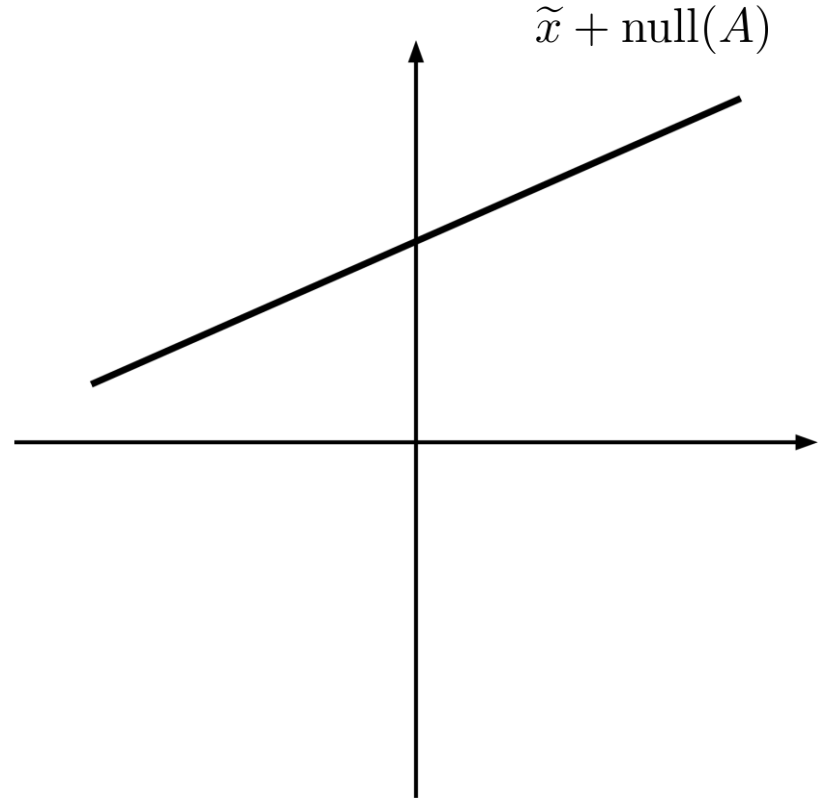
$$\hat{x} \in \arg \min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$y = Ax$ has solutions $\tilde{x} + \text{null}(A)$

$$y = A\tilde{x} \Big| \quad \Big| \{d : Ad = 0\}$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



L1-Norm Induces Sparsity

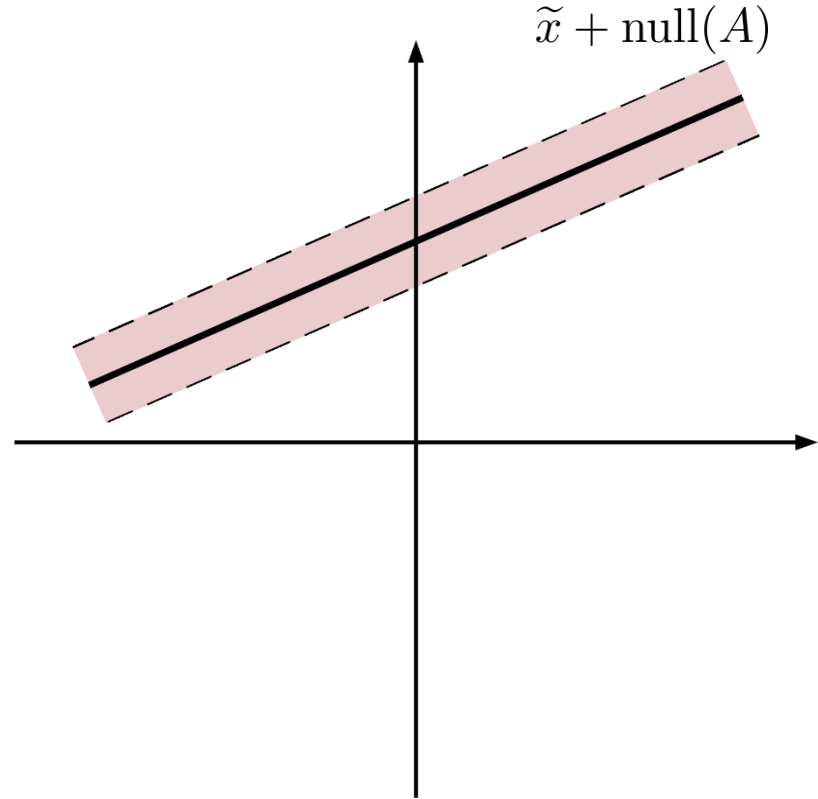
$$\hat{x} \in \arg \min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$

$$y = A\tilde{x} \Big| \Big| \{d : Ad = 0\}$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



L1-Norm Induces Sparsity

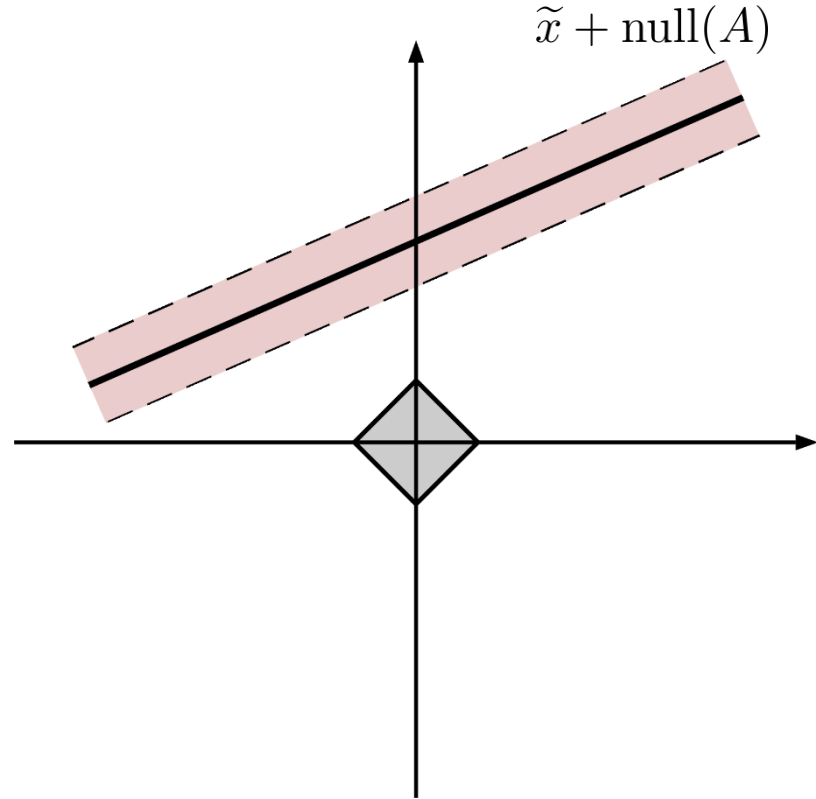
$$\hat{x} \in \arg \min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$

$$y = A\tilde{x} \quad \left| \quad \left\{ d : Ad = 0 \right\} \right.$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



L1-Norm Induces Sparsity

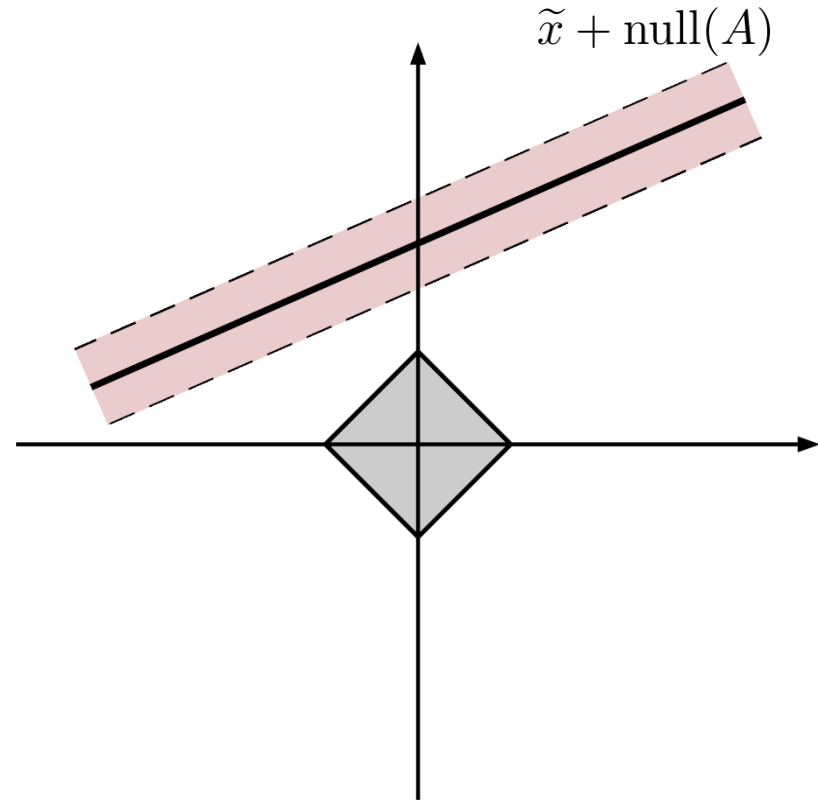
$$\hat{x} \in \arg \min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$

$$y = A\tilde{x} \quad \left| \quad \left\{ d : Ad = 0 \right\} \right.$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



L1-Norm Induces Sparsity

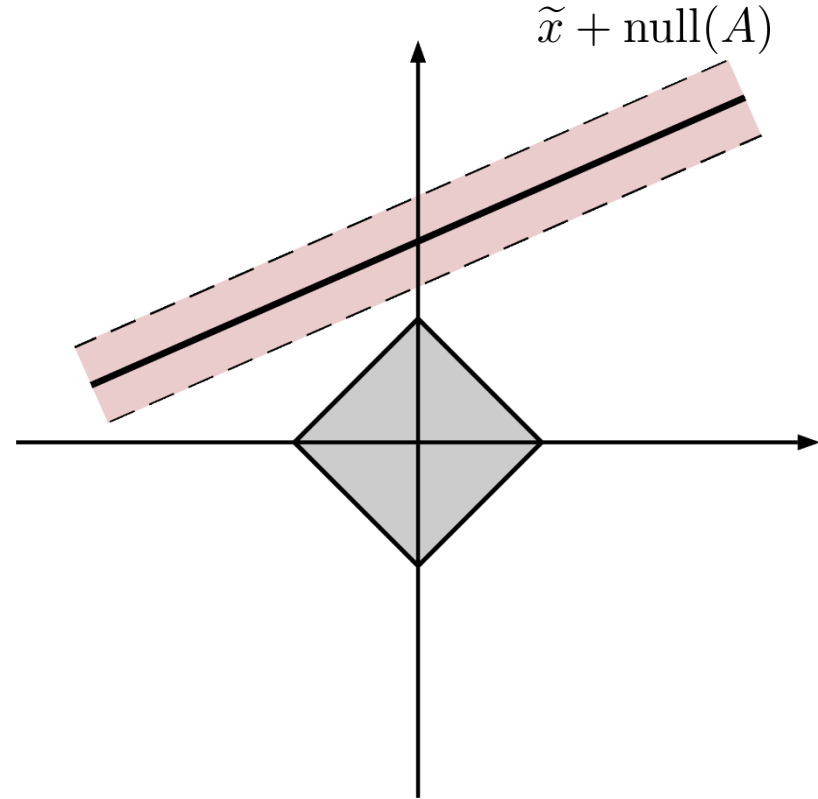
$$\hat{x} \in \arg \min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$

$$y = A\tilde{x} \quad \left| \quad \left\{ d : Ad = 0 \right\} \right.$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



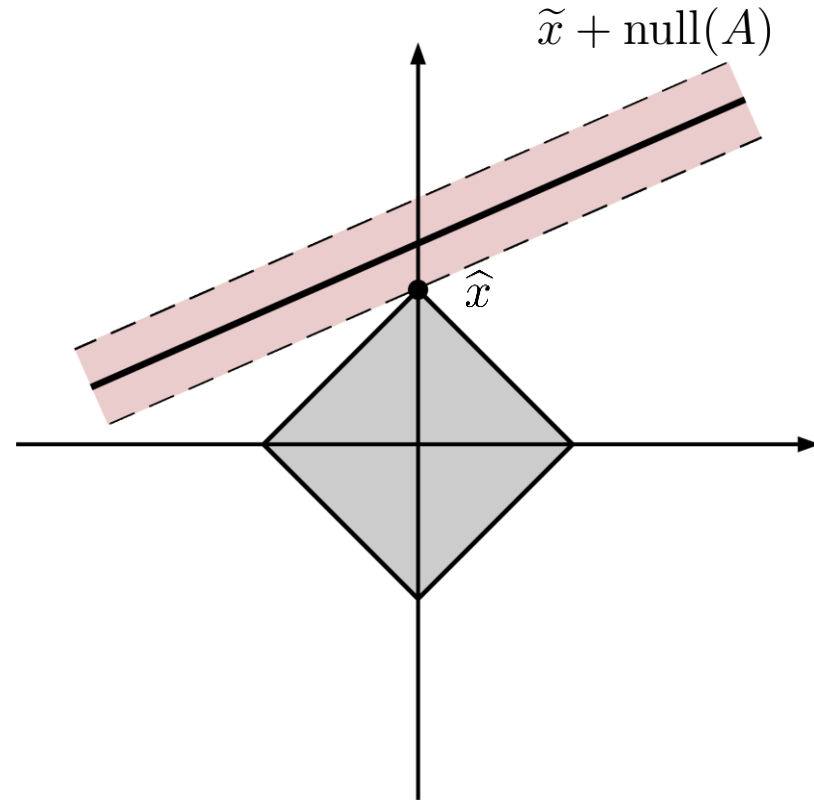
L1-Norm Induces Sparsity

$$\begin{aligned} \hat{x} \in \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & \|y - Ax\|_2 \leq \sigma \end{aligned}$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$
$$y = A\tilde{x} \quad \left| \quad \left\{ d : Ad = 0 \right\} \right.$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



L1-Norm Induces Sparsity

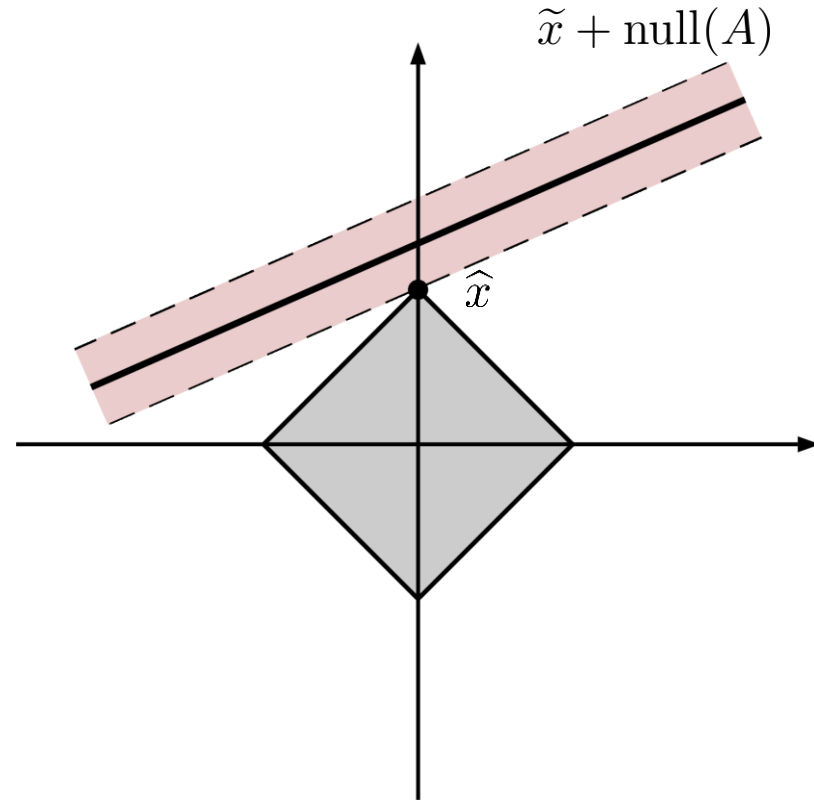
$$\hat{x} \in \arg \min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$

$$y = A\tilde{x} \quad \left| \quad \left\{ d : Ad = 0 \right\} \right.$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



What about the L2-norm?

$$\hat{x}_2 \in \arg \min_x \|x\|_2^2 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

L1-Norm Induces Sparsity

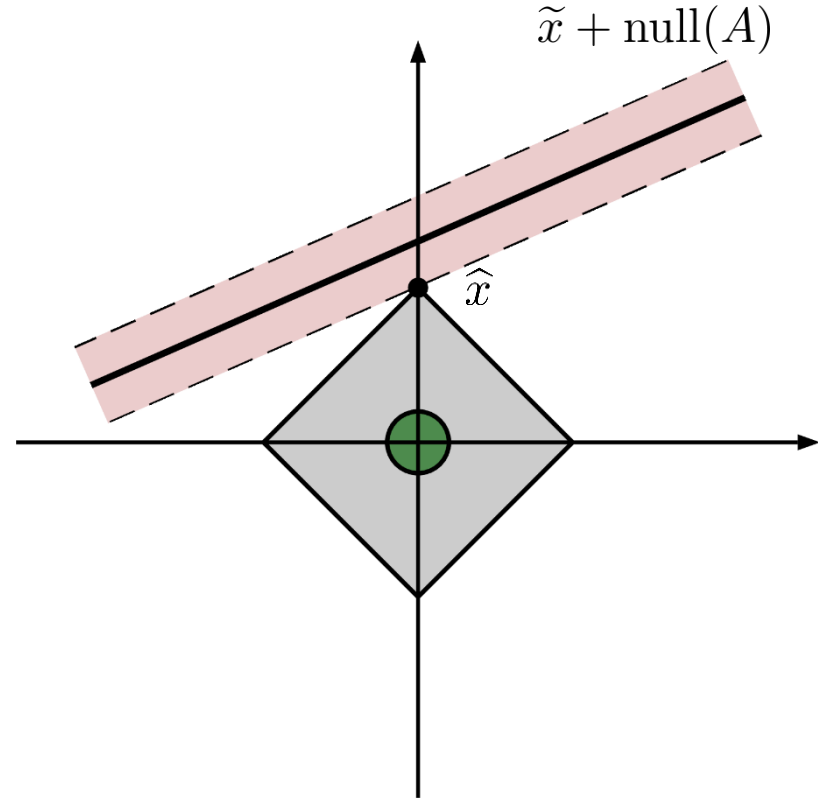
$$\hat{x} \in \arg \min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$

$$y = A\tilde{x} \quad \left| \quad \{d : Ad = 0\} \right.$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



What about the L2-norm?

$$\hat{x}_2 \in \arg \min_x \|x\|_2^2 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

L1-Norm Induces Sparsity

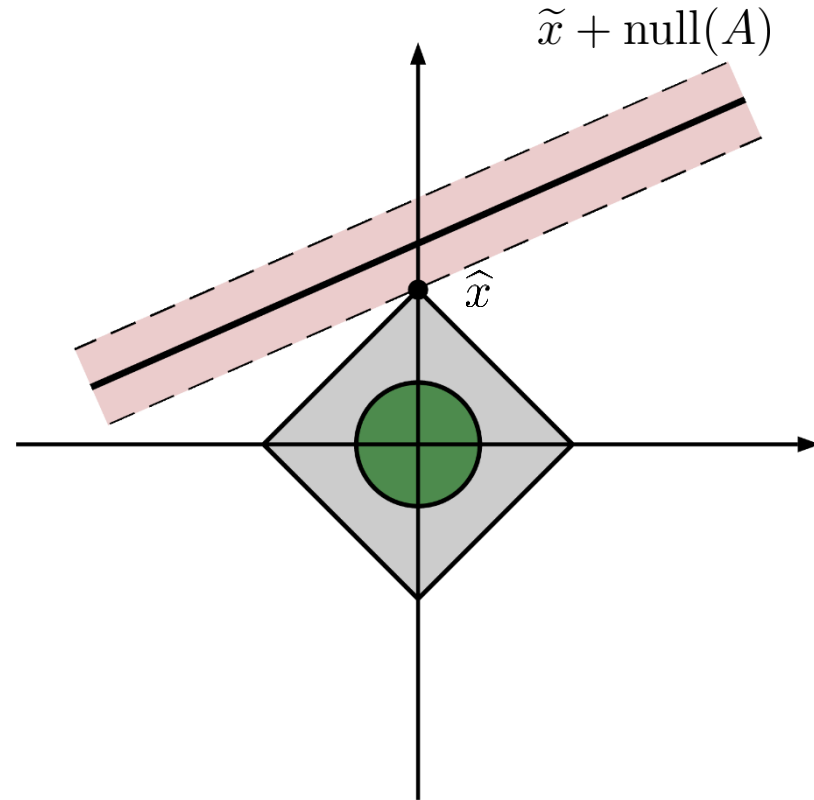
$$\hat{x} \in \arg \min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$

$$y = A\tilde{x} \quad \left| \quad \{d : Ad = 0\} \right.$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



What about the L2-norm?

$$\hat{x}_2 \in \arg \min_x \|x\|_2^2 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

L1-Norm Induces Sparsity

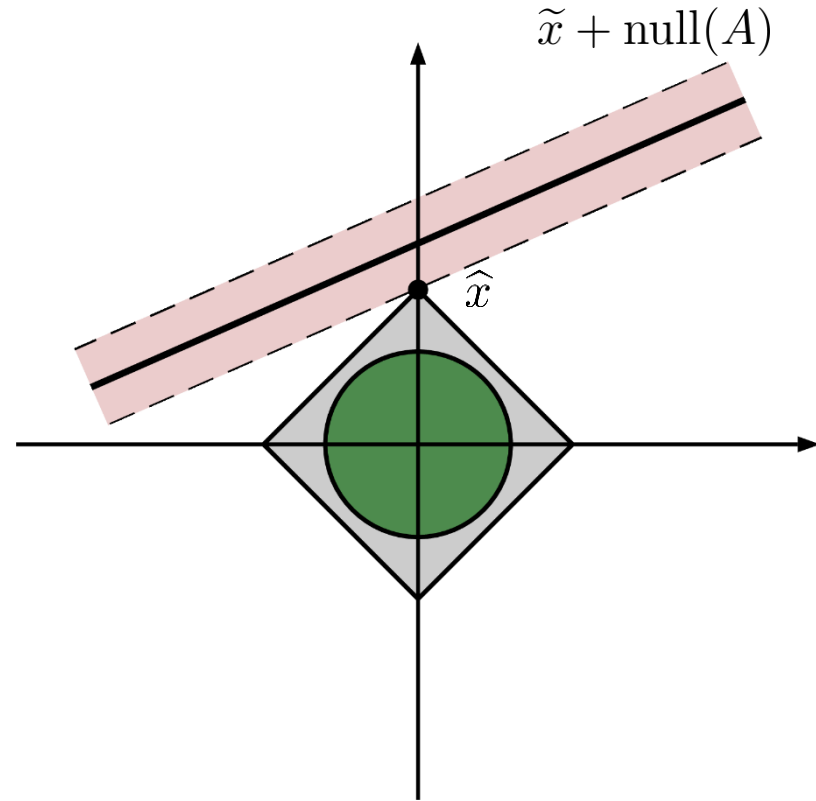
$$\hat{x} \in \arg \min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$

$$y = A\tilde{x} \quad \left| \quad \left\{ d : Ad = 0 \right\} \right.$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



What about the L2-norm?

$$\hat{x}_2 \in \arg \min_x \|x\|_2^2 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

L1-Norm Induces Sparsity

$$\hat{x} \in \arg \min_x \|x\|_1$$

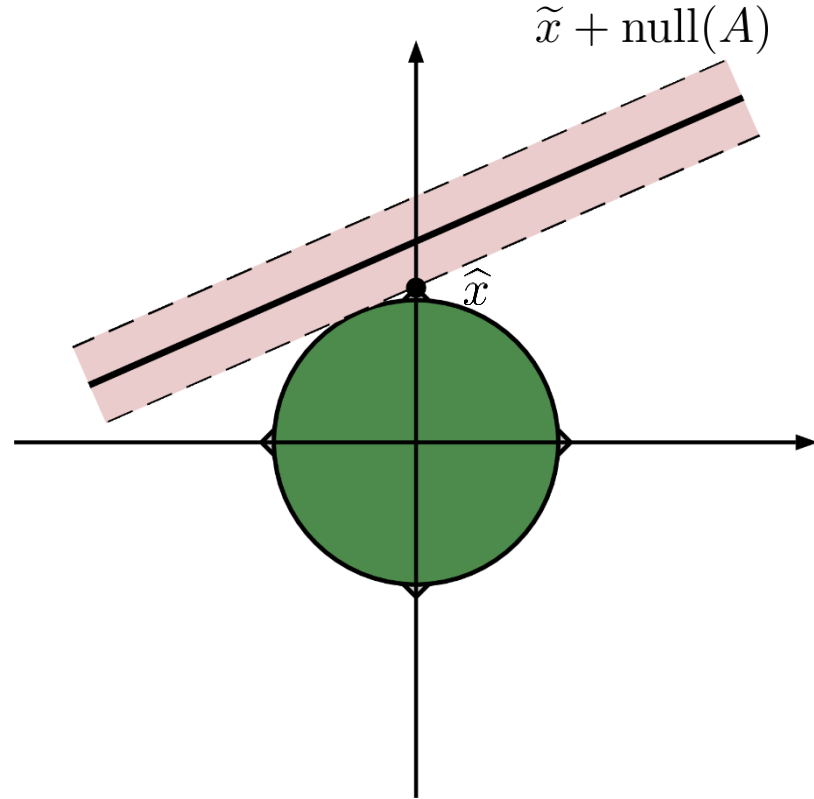
$$\text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$y = Ax$ has solutions $\tilde{x} + \text{null}(A)$

$$y = A\tilde{x} \Big| \Big| \{d : Ad = 0\}$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



What about the L2-norm?

$$\hat{x}_2 \in \arg \min_x \|x\|_2^2$$

$$\text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

L1-Norm Induces Sparsity

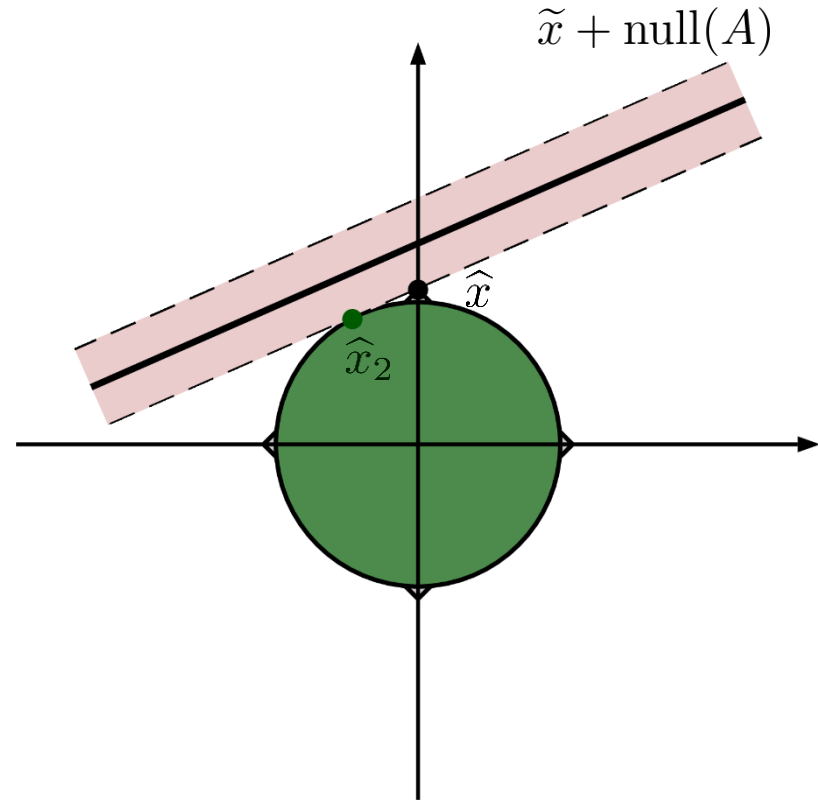
$$\hat{x} \in \arg \min_x \|x\|_1 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Assume $\sigma = 0$:

$$y = Ax \text{ has solutions } \tilde{x} + \text{null}(A)$$

$$y = A\tilde{x} \Big| \Big| \{d : Ad = 0\}$$

Assume $\sigma > 0$: *margin around* $\tilde{x} + \text{null}(A)$



What about the L2-norm?

$$\hat{x}_2 \in \arg \min_x \|x\|_2^2 \quad \text{s.t.} \quad \|y - Ax\|_2 \leq \sigma$$

Example: Compressed Sensing

Example: Compressed Sensing

$$y = A : n \times d \quad \begin{matrix} \text{sparse} \\ x^* \end{matrix}$$

Example: Compressed Sensing

$$y = A : n \times d \quad \text{sparse} \quad x^{\star}$$

Basis Pursuit

$$\begin{aligned} \hat{x} = \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & y = Ax \end{aligned}$$

Example: Compressed Sensing

Compressed Sensing (CS)

$$y = A : n \times d \quad \text{sparse} \quad x^{\star}$$

Basis Pursuit

$$\begin{aligned} \hat{x} = \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & y = Ax \end{aligned}$$

Example: Compressed Sensing

Compressed Sensing (CS)

$$y = A : n \times d \quad \begin{array}{l} \text{70-sparse} \\ \text{1000} \end{array} \quad x^*$$

iid Gaussian

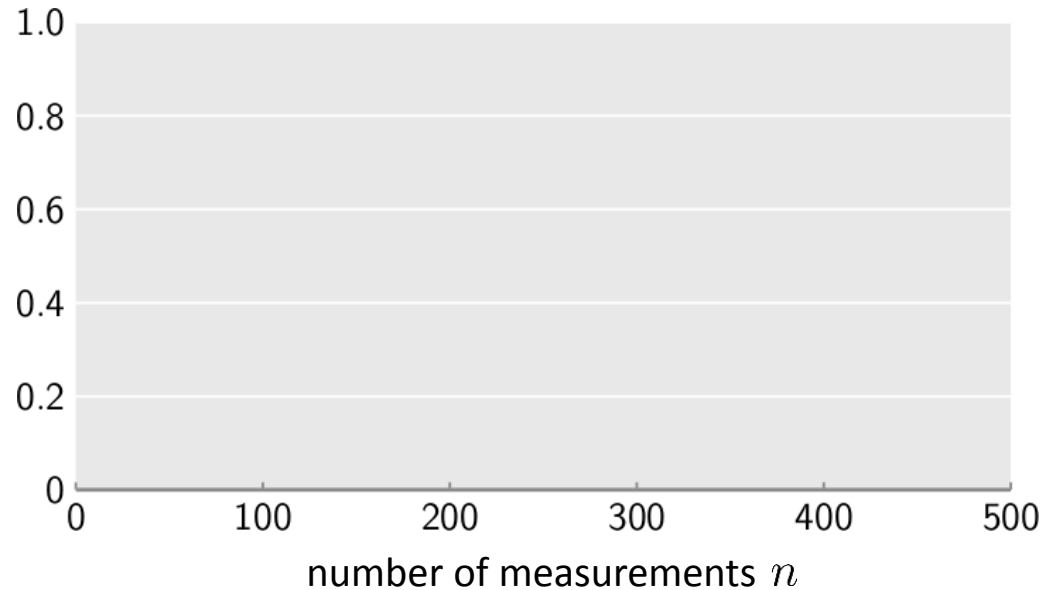
Basis Pursuit

$$\hat{x} = \arg \min_x \|x\|_1$$

s.t. $y = Ax$

Example: Compressed Sensing

Success rate (50 trials)



Compressed Sensing (CS)

$$y = A : n \times d \quad \begin{array}{l} \text{70-sparse} \\ \text{1000} \end{array} \quad x^*$$

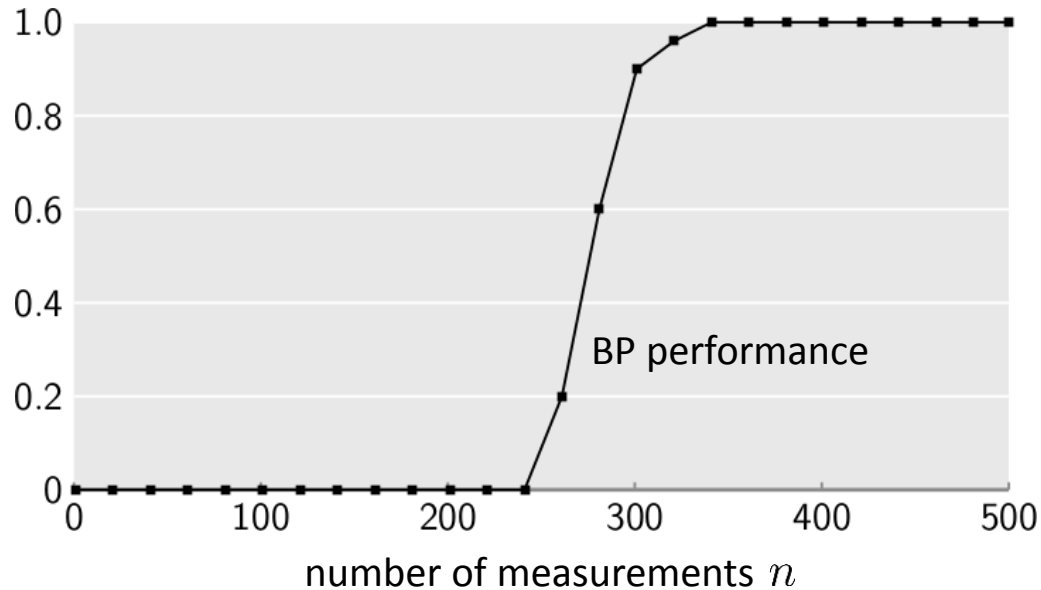
iid Gaussian

Basis Pursuit

$$\hat{x} = \arg \min_x \|x\|_1 \quad \text{s.t.} \quad y = Ax$$

Example: Compressed Sensing

Success rate (50 trials)



Compressed Sensing (CS)

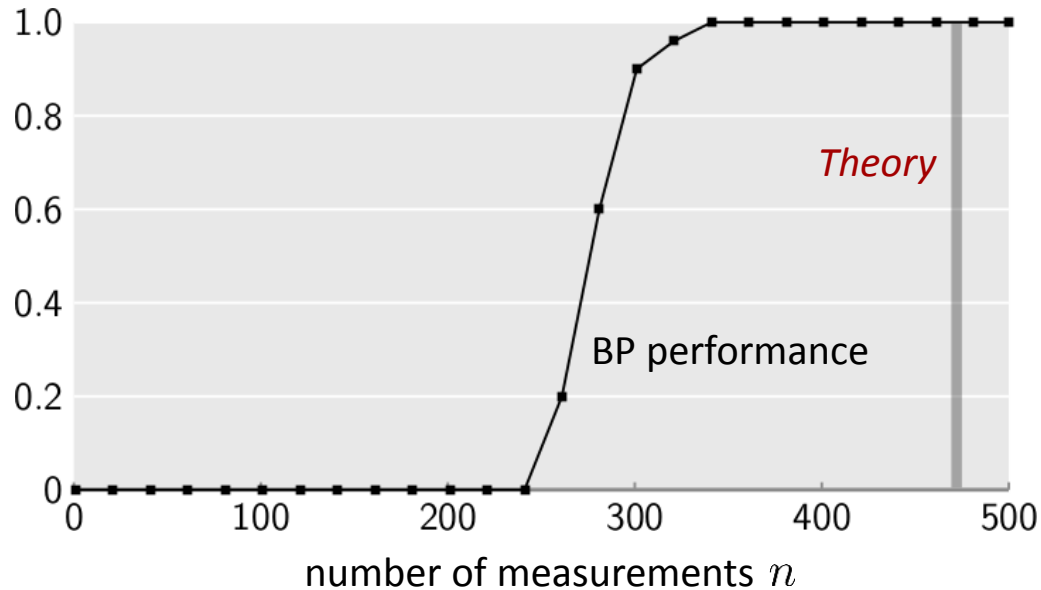
$$y = A : n \times d \quad \begin{array}{l} \text{70-sparse} \\ \text{iid Gaussian} \\ 1000 \end{array} \quad x^*$$

Basis Pursuit

$$\hat{x} = \arg \min_x \|x\|_1 \quad \text{s.t.} \quad y = Ax$$

Example: Compressed Sensing

Success rate (50 trials)



Compressed Sensing (CS)

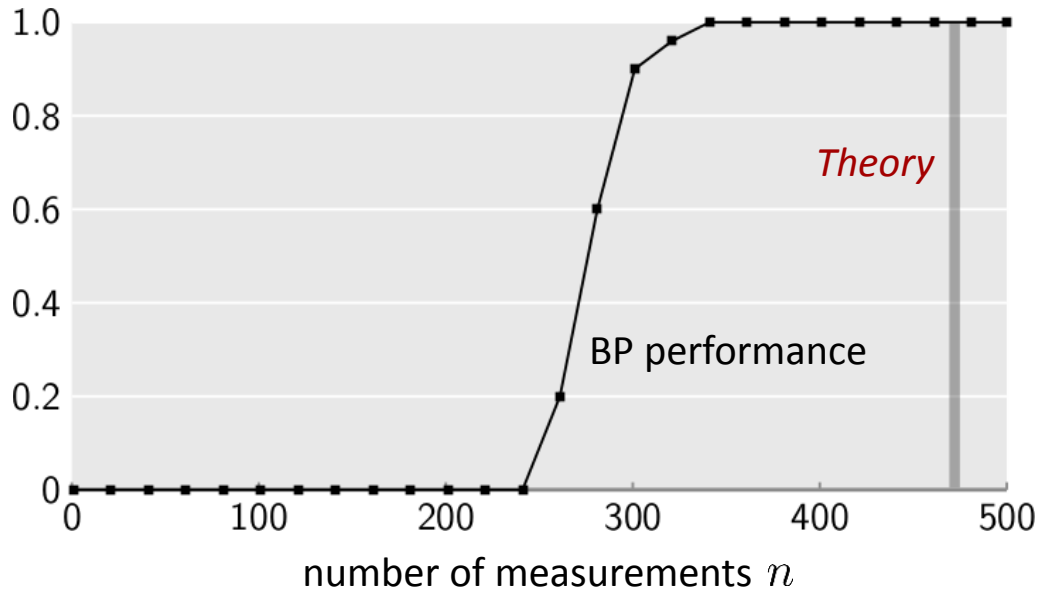
$$y = A : n \times d \quad \begin{array}{l} \text{70-sparse} \\ \text{iid Gaussian} \\ 1000 \end{array} \quad x^*$$

Basis Pursuit

$$\hat{x} = \arg \min_x \|x\|_1 \quad \text{s.t.} \quad y = Ax$$

Example: Compressed Sensing

Success rate (50 trials)



Compressed Sensing (CS)

$$y = A : n \times d \quad \begin{array}{l} \text{70-sparse} \\ \text{iid Gaussian} \\ 1000 \end{array} \quad x^*$$

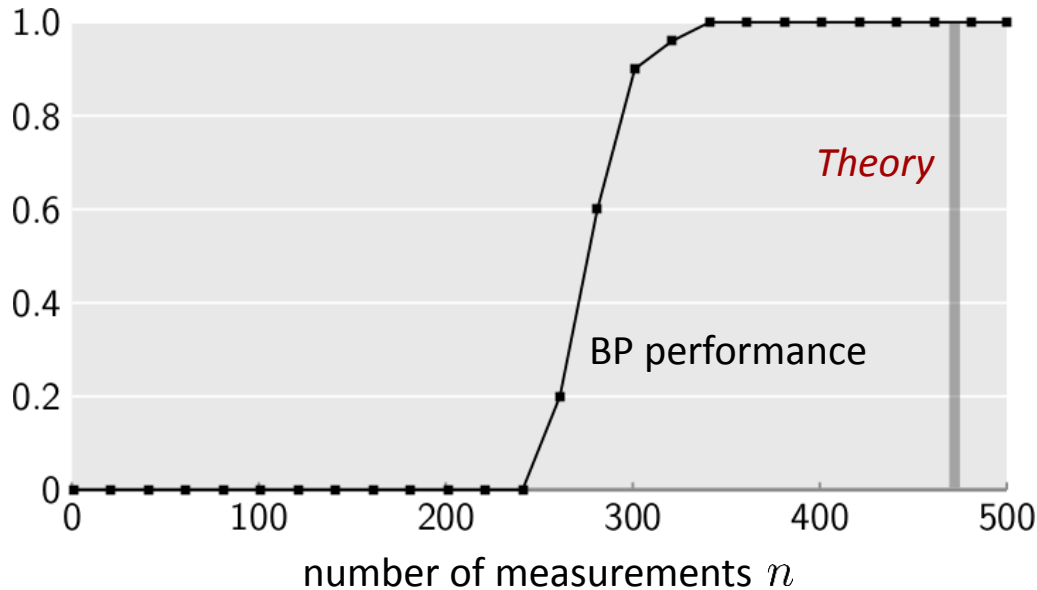
Basis Pursuit

$$\hat{x} = \arg \min_x \|x\|_1 \quad \text{s.t.} \quad y = Ax$$

Theorem [Chandrasekaran et al. 12']

Example: Compressed Sensing

Success rate (50 trials)



Compressed Sensing (CS)

$$y = A : n \times d \quad \begin{array}{l} \text{70-sparse} \\ \text{iid Gaussian} \\ 1000 \end{array} \quad x^*$$

Basis Pursuit

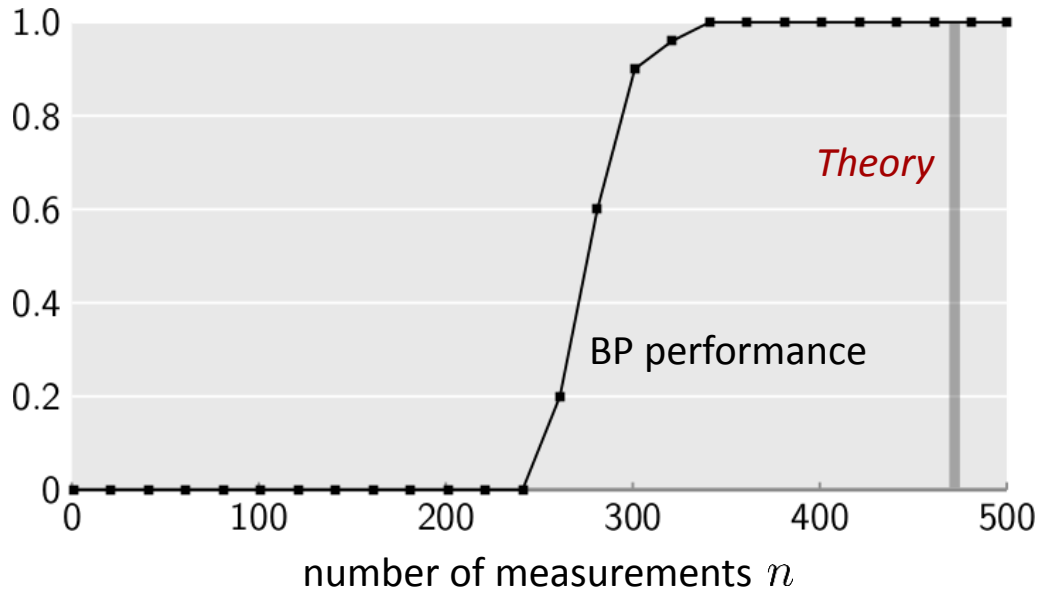
$$\hat{x} = \arg \min_x \|x\|_1 \quad \text{s.t.} \quad y = Ax$$

Theorem [Chandrasekaran et al. 12']

$$x^* \in \mathbb{R}^d$$

Example: Compressed Sensing

Success rate (50 trials)



Compressed Sensing (CS)

$$y = A : n \times d \quad \begin{array}{l} \text{70-sparse} \\ \text{iid Gaussian} \\ 1000 \end{array} \quad x^*$$

Basis Pursuit

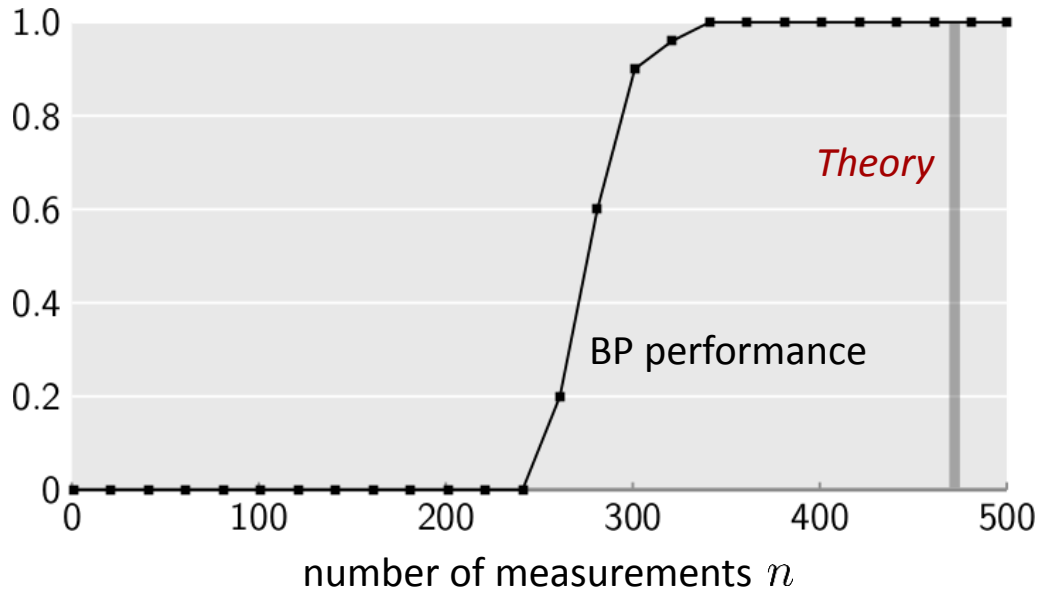
$$\hat{x} = \arg \min_x \|x\|_1 \quad \text{s.t.} \quad y = Ax$$

Theorem [Chandrasekaran et al. 12']

$x^* \in \mathbb{R}^d$ *unknown*, but s -sparse

Example: Compressed Sensing

Success rate (50 trials)



Compressed Sensing (CS)

$$y = A : n \times d$$

iid Gaussian

1000

70-sparse

x^*

Basis Pursuit

$$\hat{x} = \arg \min_x \|x\|_1 \quad \text{s.t.} \quad y = Ax$$

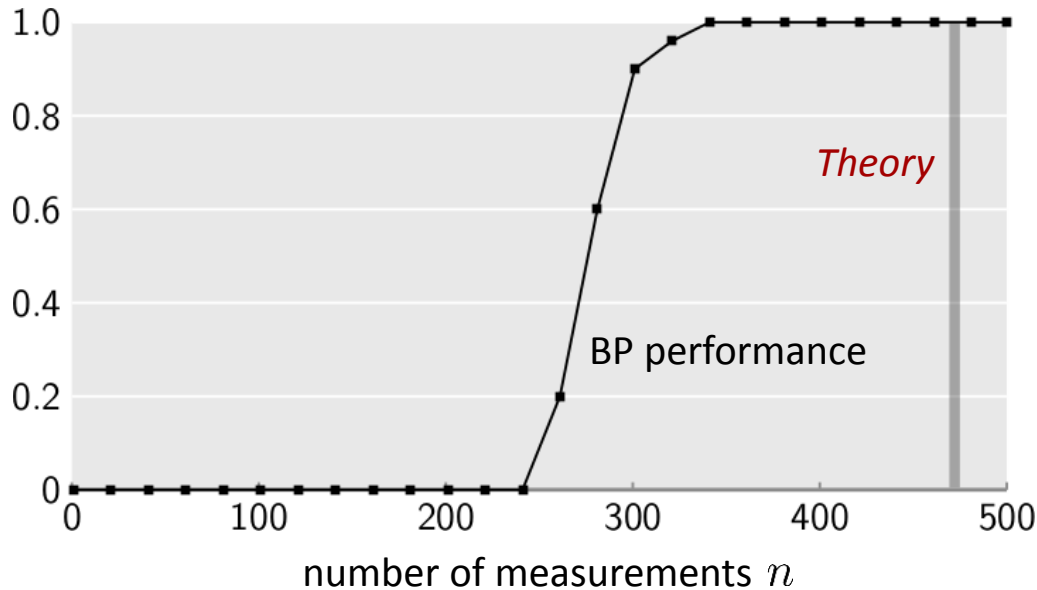
Theorem [Chandrasekaran et al. 12']

$x^* \in \mathbb{R}^d$ *unknown*, but s -sparse

$y = Ax^*$ *measurements*

Example: Compressed Sensing

Success rate (50 trials)



Compressed Sensing (CS)

$$y = A : n \times d \quad \begin{array}{l} \text{70-sparse} \\ \text{iid Gaussian} \\ 1000 \end{array} \quad x^*$$

Basis Pursuit

$$\hat{x} = \arg \min_x \|x\|_1 \quad \text{s.t.} \quad y = Ax$$

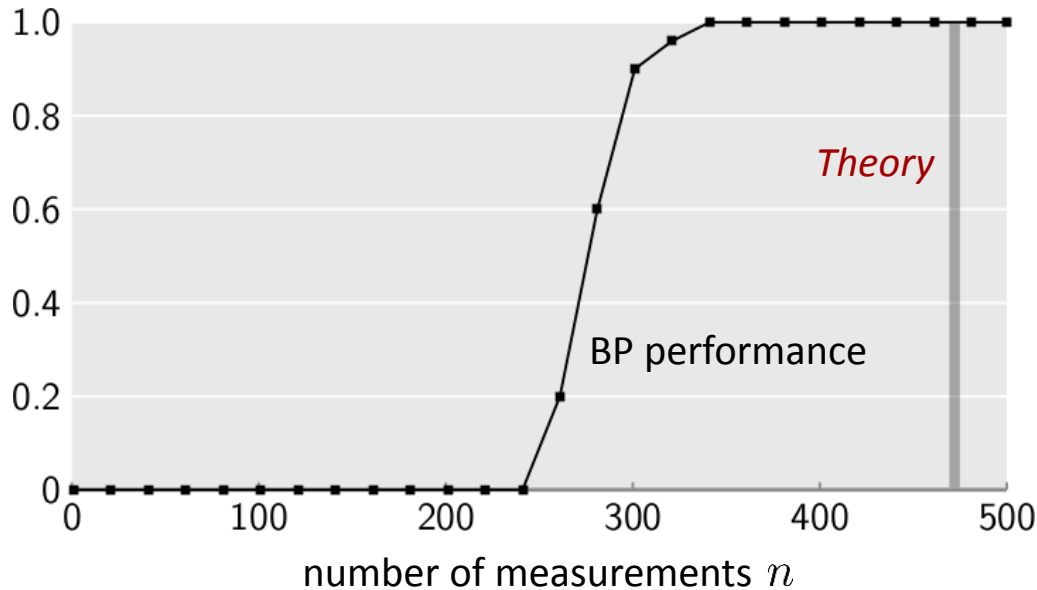
Theorem [Chandrasekaran et al. 12']

$x^* \in \mathbb{R}^d$ *unknown*, but s -sparse

$$y = Ax^* \quad \begin{array}{l} \text{iid entries } \mathcal{N}(0, 1/n) \\ \text{measurements} \end{array}$$

Example: Compressed Sensing

Success rate (50 trials)



Compressed Sensing (CS)

$$y = A : n \times d \quad \begin{array}{l} \text{70-sparse} \\ \text{iid Gaussian} \\ 1000 \end{array} \quad x^*$$

Basis Pursuit

$$\hat{x} = \arg \min_x \|x\|_1 \quad \text{s.t.} \quad y = Ax$$

Theorem [Chandrasekaran et al. 12']

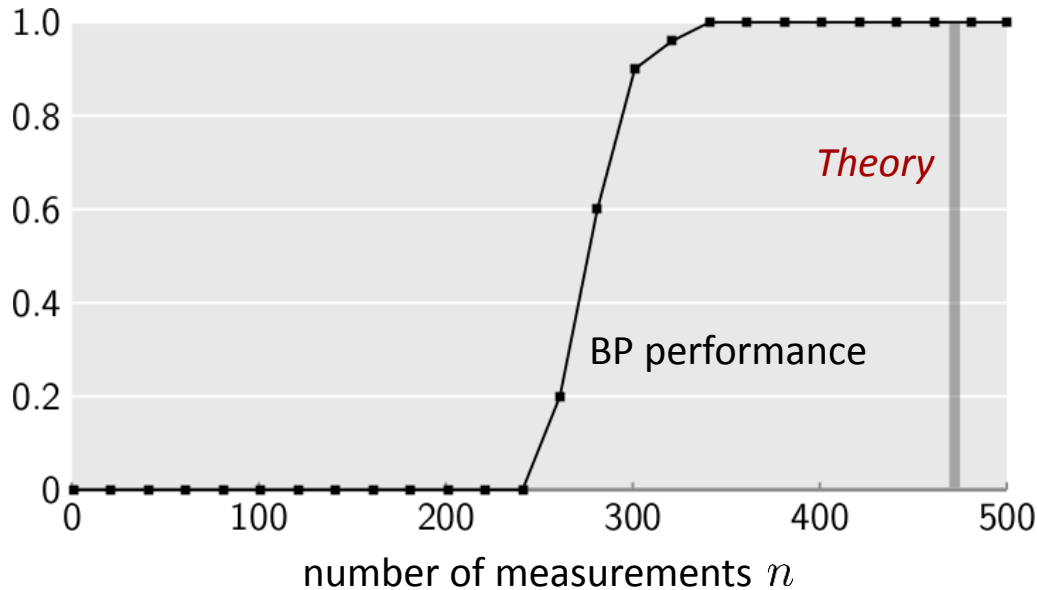
$x^* \in \mathbb{R}^d$ **unknown**, but s -sparse

$$y = Ax^* \quad \begin{array}{l} \text{iid entries } \mathcal{N}(0, 1/n) \\ \text{measurements} \end{array}$$

$$n \geq 2s \log\left(\frac{d}{s}\right) + \frac{7}{5}s + 1$$

Example: Compressed Sensing

Success rate (50 trials)



Compressed Sensing (CS)

$$y = A : n \times d \quad \begin{array}{l} \text{70-sparse} \\ \text{iid Gaussian} \\ 1000 \end{array} \quad x^*$$

Basis Pursuit

$$\hat{x} = \arg \min_x \|x\|_1 \quad \text{s.t.} \quad y = Ax$$

Theorem [Chandrasekaran et al. 12']

$x^* \in \mathbb{R}^d$ *unknown*, but s -sparse

$$y = Ax^* \quad \begin{array}{l} \text{iid entries } \mathcal{N}(0, 1/n) \\ \text{measurements} \end{array}$$

$$n \geq 2s \log\left(\frac{d}{s}\right) + \frac{7}{5}s + 1$$

\implies

$$x^* = \arg \min_x \|x\|_1 \quad \text{w.h.p.} \quad \text{s.t.} \quad y = Ax$$

Application: Image Reconstruction

Application: Image Reconstruction



Application: Image Reconstruction



256×496

Application: Image Reconstruction



$$256 \times 496 \quad \Rightarrow \quad z^* \in \mathbb{R}^{126976}$$

Application: Image Reconstruction



$$256 \times 496 \quad \Rightarrow \quad z^* \in \mathbb{R}^{126976}$$

not sparse

Application: Image Reconstruction



$$256 \times 496 \quad \Rightarrow \quad z^* \in \mathbb{R}^{126976}$$

not sparse

Natural images have sparse representations

Application: Image Reconstruction



$$256 \times 496 \quad \Rightarrow \quad z^* \in \mathbb{R}^{126976}$$

not sparse

Natural images have sparse representations

$$z^* = \Psi x^*$$

Application: Image Reconstruction



$$256 \times 496 \quad \Rightarrow \quad z^* \in \mathbb{R}^{126976}$$

not sparse

Natural images have sparse representations

$$z^* = \Psi x^*$$

| *sparse* or near-sparse

Application: Image Reconstruction



$$256 \times 496 \quad \Rightarrow \quad z^* \in \mathbb{R}^{126976}$$

not sparse

Natural images have sparse representations

$$z^* = \Psi x^*$$

sparse or near-sparse
dictionary (wavelet, DCT, gradient space)

Application: Image Reconstruction



$$256 \times 496 \quad \Rightarrow \quad z^* \in \mathbb{R}^{126976}$$

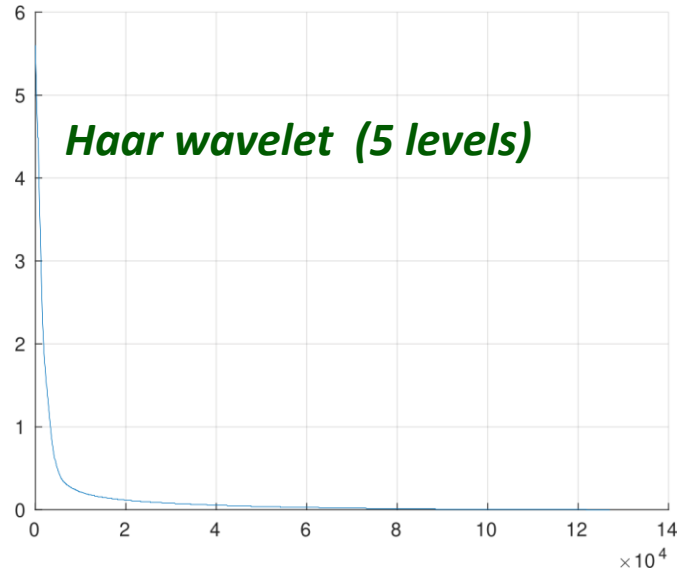
not sparse

Natural images have sparse representations

$$z^* = \Psi x^*$$

sparse or near-sparse
dictionary (wavelet, DCT, gradient space)

Ordered coefficients: $|x_i^*|$

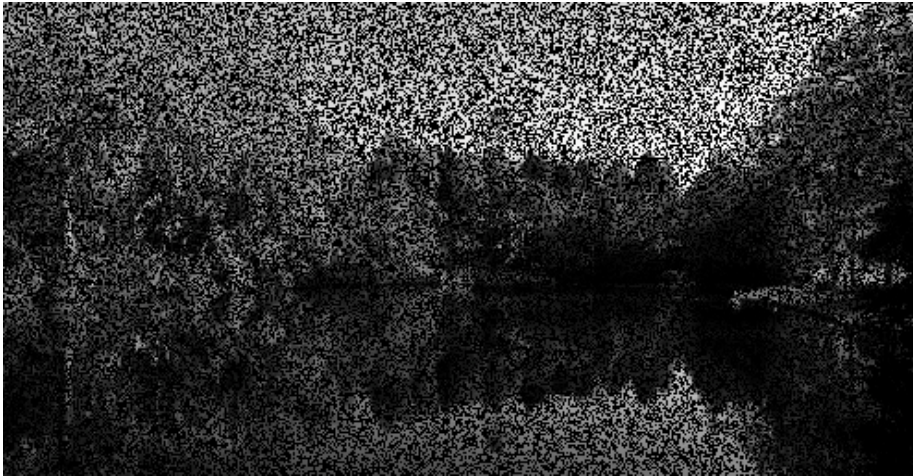


Application: Image Reconstruction

Suppose we observe *only 50%* of pixels

Application: Image Reconstruction

Suppose we observe *only 50%* of pixels



y

Application: Image Reconstruction

Suppose we observe *only 50%* of pixels

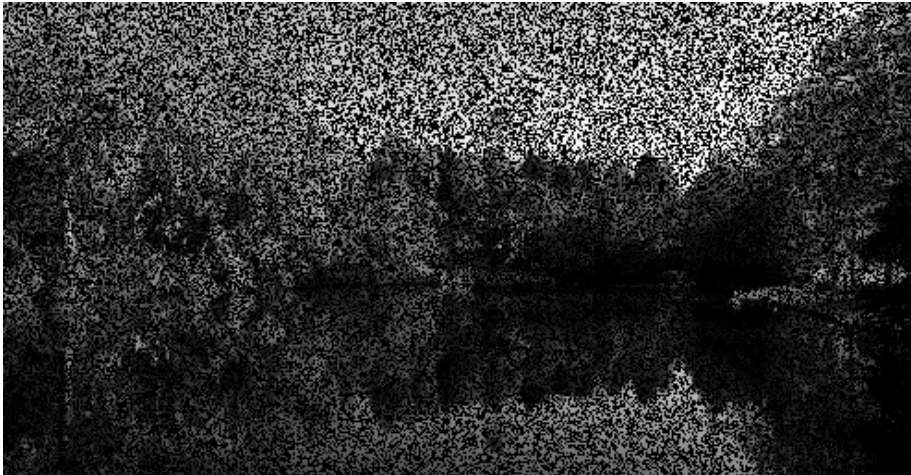
 y

Solve

$$\begin{aligned} \hat{x} = & \arg \min_x \|x\|_1 \\ \text{s.t.} \quad & y = \Phi \Psi x \end{aligned}$$

Application: Image Reconstruction

Suppose we observe *only 50%* of pixels

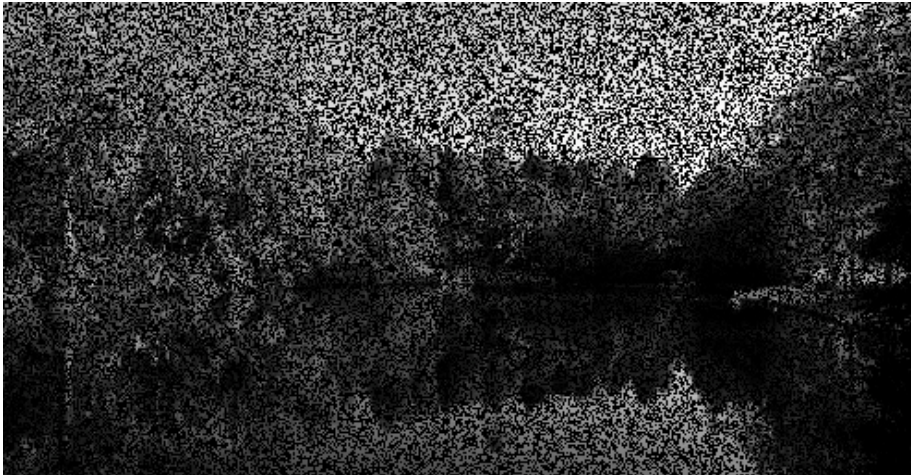
 y

Solve

$$\begin{aligned} \hat{x} = & \arg \min_x \|x\|_1 \\ \text{s.t.} \quad & y = \Phi \Psi x \\ & \quad \quad \quad | \text{wavelet} \end{aligned}$$

Application: Image Reconstruction

Suppose we observe *only 50%* of pixels

 y

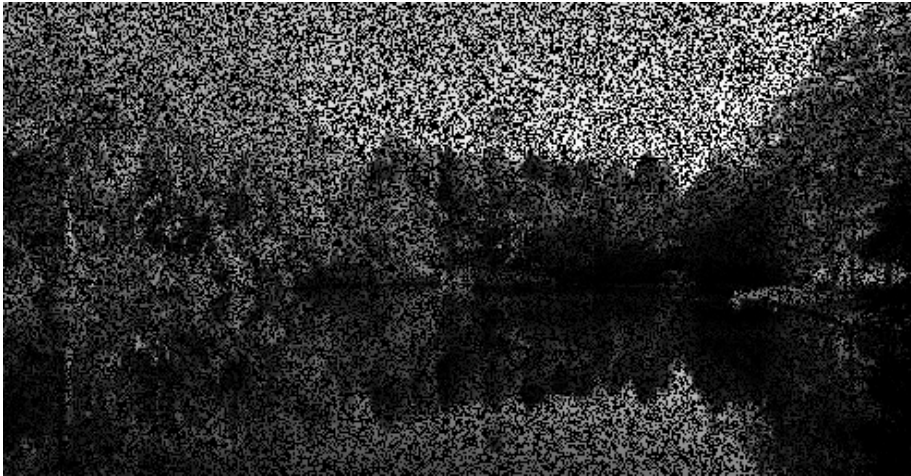
Solve

$$\begin{aligned} \hat{x} = \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & y = \Phi \Psi x \end{aligned}$$

| wavelet
| observed indices

Application: Image Reconstruction

Suppose we observe *only 50%* of pixels

 y

Solve

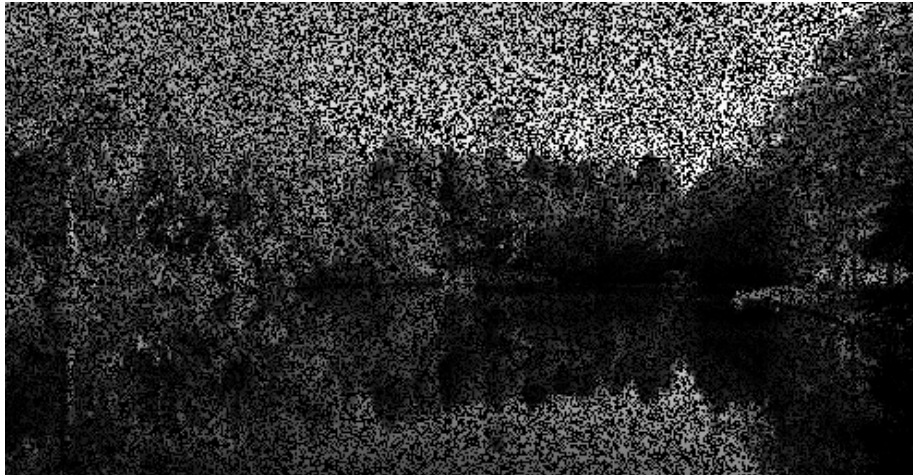
$$\begin{aligned} \hat{x} = \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & y = \Phi \Psi x \end{aligned}$$

| wavelet
| observed indices

 \hat{x}

Application: Image Reconstruction

Suppose we observe *only 50%* of pixels

 y

Solve

$$\begin{aligned} \hat{x} = \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & y = \Phi \Psi x \end{aligned}$$

| wavelet
| observed indices

 \hat{x}

PSNR: 21.31 dB

Application: Image Reconstruction

Solve

$$\begin{aligned} \hat{x} = \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & y = \Phi \Psi x \end{aligned}$$

| wavelet
partial DFT

Application: Image Reconstruction

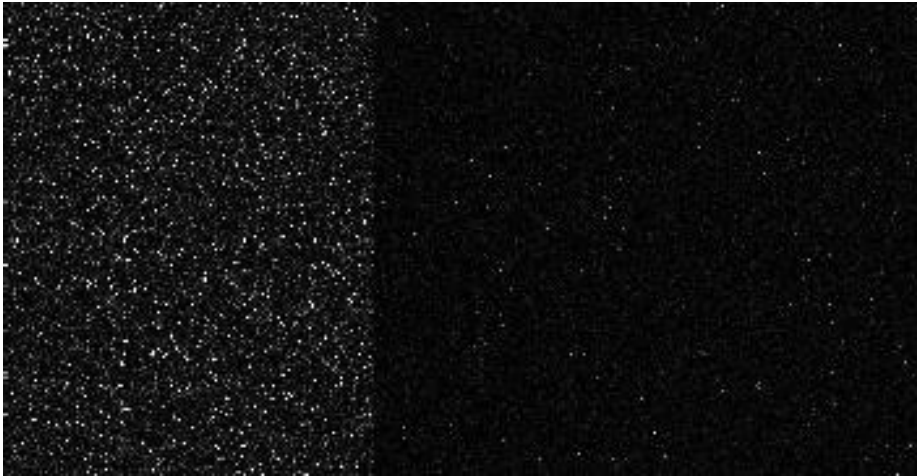
Solve

$$\begin{aligned} \hat{x} = \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & y = \Phi \Psi x \end{aligned}$$

| wavelet
| *partial DFT*

each entry of y has info from entire image

Application: Image Reconstruction

 y

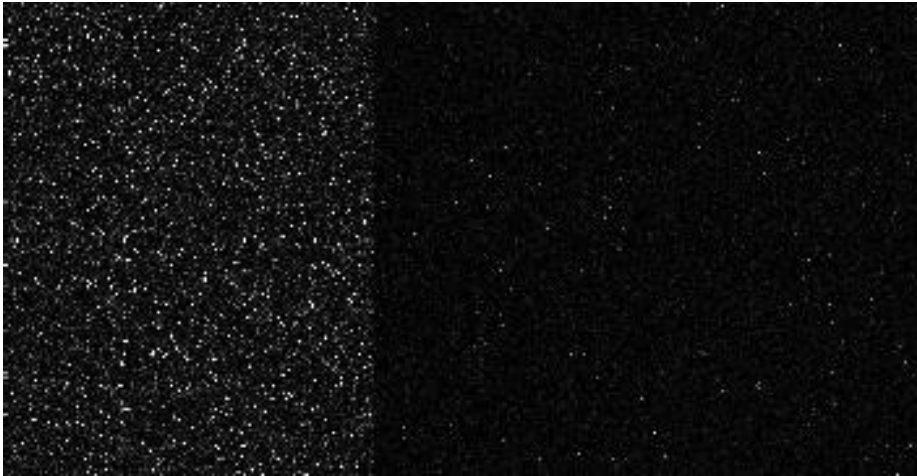
Solve

$$\begin{aligned} \hat{x} = \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & y = \Phi \Psi x \end{aligned}$$

| wavelet
| *partial DFT*

each entry of y has info from entire image

Application: Image Reconstruction

 y

Solve

$$\begin{aligned} \hat{x} = \arg \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & y = \Phi \Psi x \end{aligned}$$

| wavelet
| *partial DFT*

each entry of y has info from entire image

 \hat{x}

PSNR: 24.93 dB

Outline

- *Motivation: Hypothesis Testing in High-Dimensions*
- *Introduction to LASSO and other sparsity problems*
- Gaussian graphical model selection
- Matrix completion

Gaussian Graphical Model Selection

Gaussian Graphical Model Selection

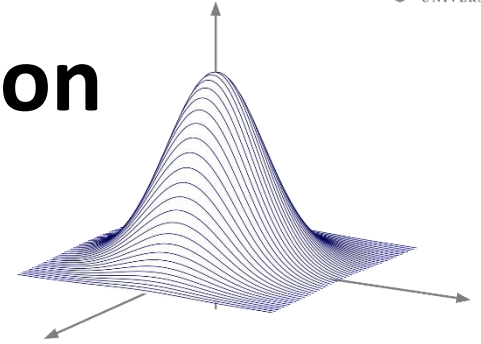
$$X = (X_1, X_2, \dots, X_d)$$

Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

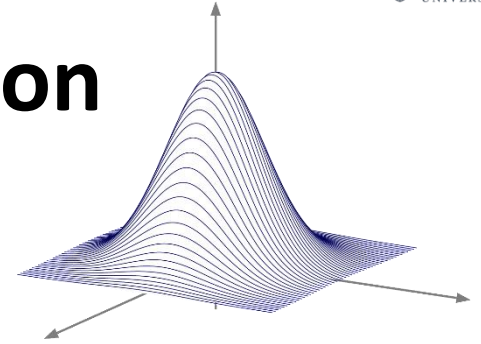
Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$



Gaussian Graphical Model Selection

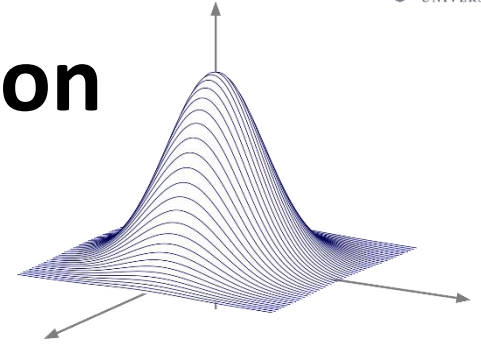
$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$



Problem: given n iid observations of X , denoted $X^{(1)}, \dots, X^{(n)}$, estimate Σ^*

Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

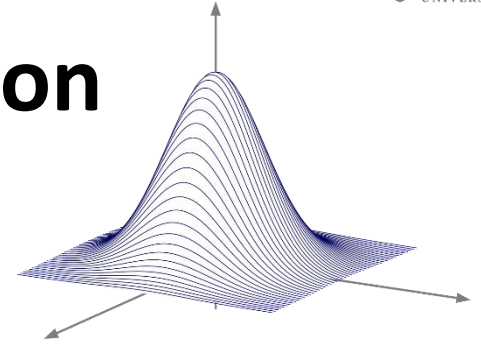


Problem: given n iid observations of X , denoted $X^{(1)}, \dots, X^{(n)}$, estimate Σ^*

Assumption:

Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

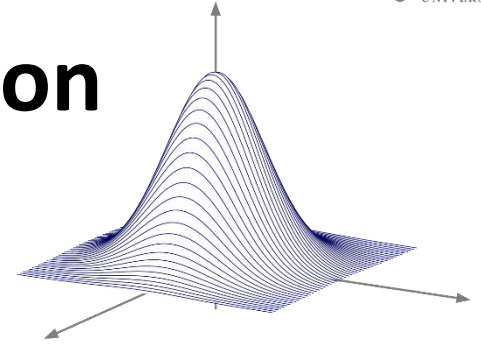


Problem: given n iid observations of X , denoted $X^{(1)}, \dots, X^{(n)}$, estimate Σ^*

Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent

Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$



Problem: given n iid observations of X , denoted $X^{(1)}, \dots, X^{(n)}$, estimate Σ^*

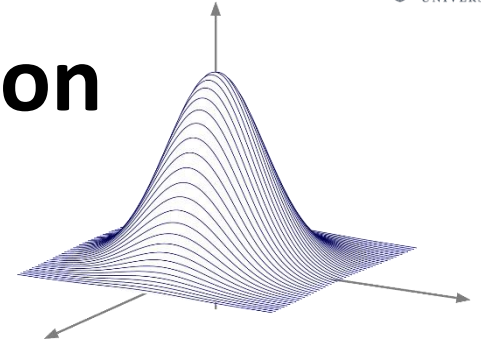
Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent



precision matrix $\Theta^* := (\Sigma^*)^{-1}$ is *sparse*

Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

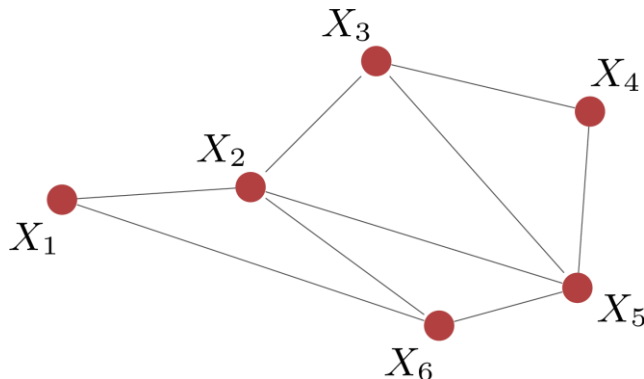


Problem: given n iid observations of X , denoted $X^{(1)}, \dots, X^{(n)}$, estimate Σ^*

Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent

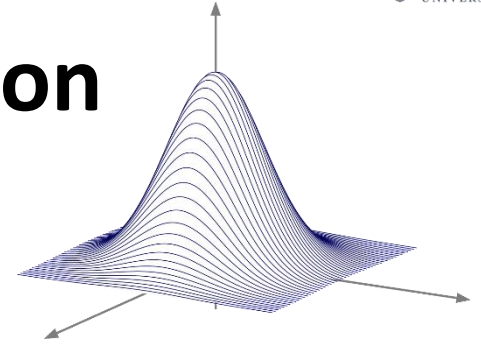


precision matrix $\Theta^* := (\Sigma^*)^{-1}$ is *sparse*



Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

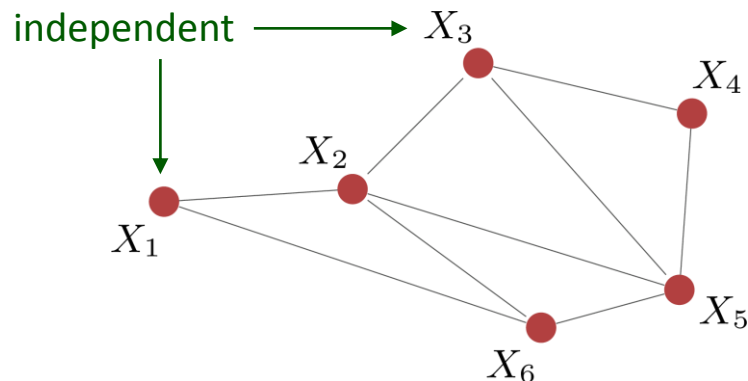


Problem: given n iid observations of X , denoted $X^{(1)}, \dots, X^{(n)}$, estimate Σ^*

Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent

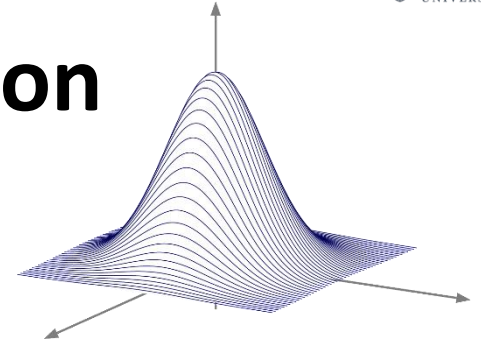


precision matrix $\Theta^* := (\Sigma^*)^{-1}$ is *sparse*



Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

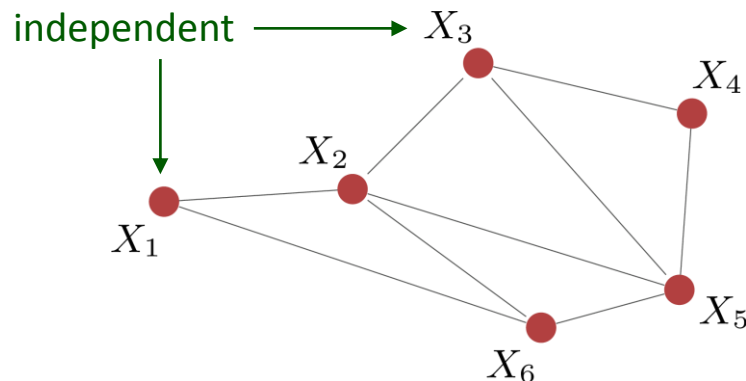


Problem: given n iid observations of X , denoted $X^{(1)}, \dots, X^{(n)}$, estimate Σ^*

Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent

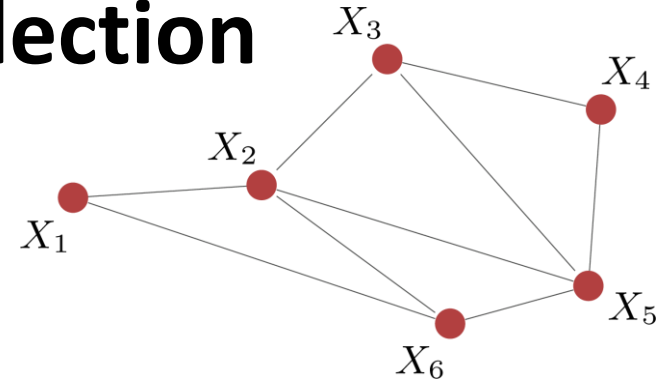


precision matrix $\Theta^* := (\Sigma^*)^{-1}$ is *sparse*



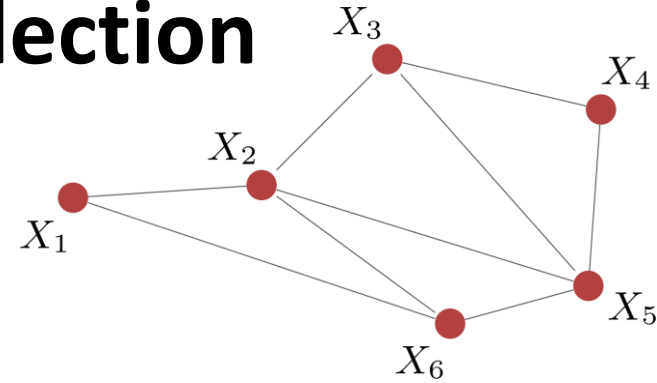
$$\Theta^* = \begin{bmatrix} \star & \star & 0 & 0 & 0 & \star \\ \star & \star & \star & 0 & \star & \star \\ 0 & \star & \star & \star & \star & 0 \\ 0 & 0 & \star & \star & \star & 0 \\ 0 & \star & \star & \star & \star & \star \\ \star & \star & 0 & 0 & \star & \star \end{bmatrix}$$

Gaussian Graphical Model Selection



Gaussian Graphical Model Selection

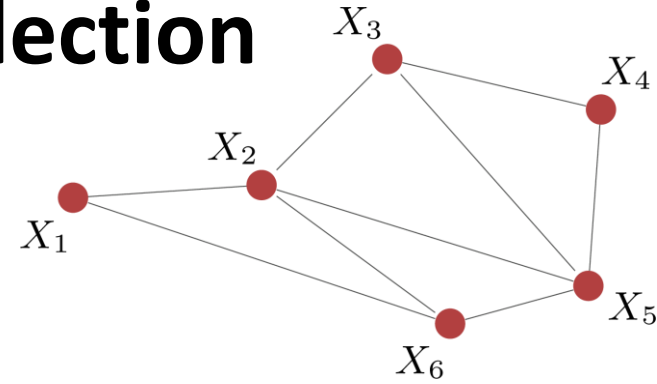
$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$



Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

$$\text{pdf: } f_X(x; \Theta^*) = \frac{\sqrt{\det \Theta^*}}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{1}{2}x^\top \Theta^* x\right)$$

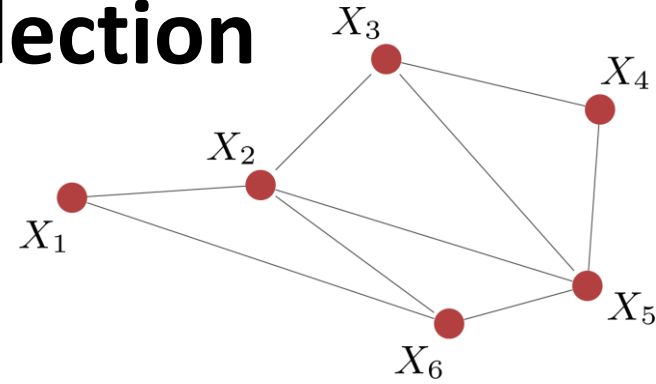


Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

$$\text{pdf: } f_X(x; \Theta^*) = \frac{\sqrt{\det \Theta^*}}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{1}{2}x^\top \Theta^* x\right)$$

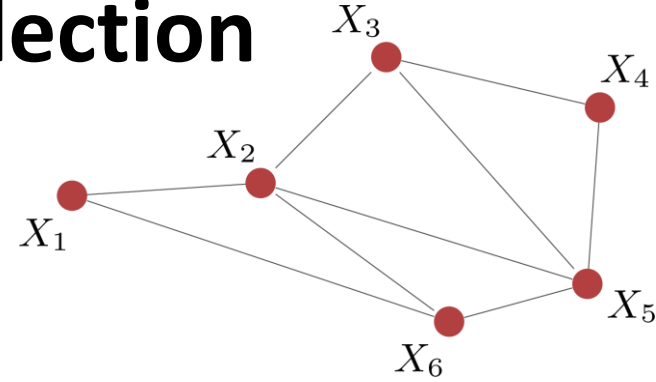
*Maximum likelihood estimator of Θ^**



Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

$$\text{pdf: } f_X(x; \Theta^*) = \frac{\sqrt{\det \Theta^*}}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{1}{2}x^\top \Theta^* x\right)$$



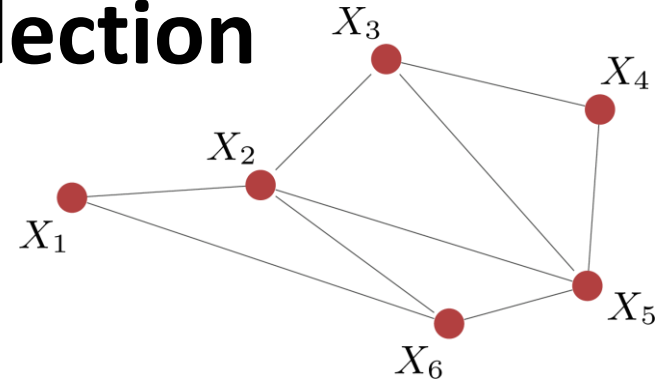
*Maximum likelihood estimator of Θ^**

$$\hat{\Theta}_{\text{ML}} = \arg \max_{\Theta} \log \prod_{i=1}^n f(x^{(1)}, \dots, x^{(n)}; \Theta)$$

Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

$$\text{pdf: } f_X(x; \Theta^*) = \frac{\sqrt{\det \Theta^*}}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{1}{2}x^\top \Theta^* x\right)$$



*Maximum likelihood estimator of Θ^**

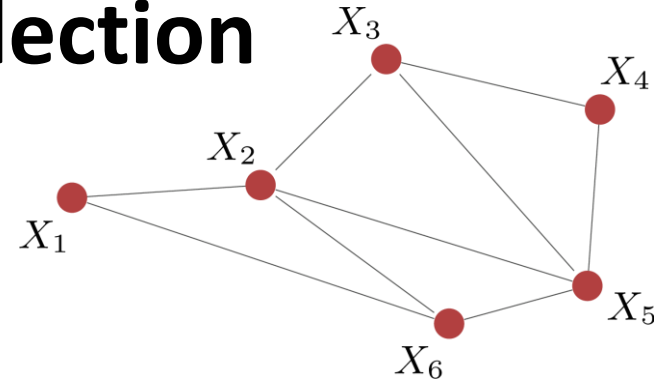
$$\hat{\Theta}_{\text{ML}} = \arg \max_{\Theta} \log \prod_{i=1}^n f(x^{(1)}, \dots, x^{(n)}; \Theta)$$

$$= \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n)$$

Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

$$\text{pdf: } f_X(x; \Theta^*) = \frac{\sqrt{\det \Theta^*}}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{1}{2}x^\top \Theta^* x\right)$$



*Maximum likelihood estimator of Θ^**

$$\hat{\Theta}_{\text{ML}} = \arg \max_{\Theta} \log \prod_{i=1}^n f(x^{(1)}, \dots, x^{(n)}; \Theta)$$

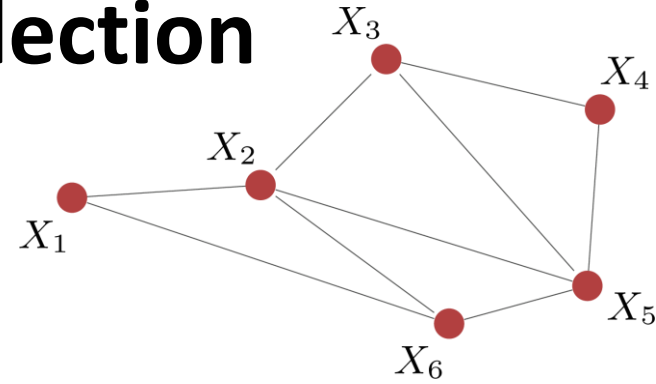
$$= \arg \min_{\Theta} -\log \det \Theta + \text{tr}\left(\Theta \hat{\Sigma}_n\right)$$

sample covariance matrix

Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

$$\text{pdf: } f_X(x; \Theta^*) = \frac{\sqrt{\det \Theta^*}}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{1}{2}x^\top \Theta^* x\right)$$



*Maximum likelihood estimator of Θ^**

$$\hat{\Theta}_{\text{ML}} = \arg \max_{\Theta} \log \prod_{i=1}^n f(x^{(1)}, \dots, x^{(n)}; \Theta)$$

$$= \arg \min_{\Theta} -\log \det \Theta + \text{tr}\left(\Theta \hat{\Sigma}_n\right)$$

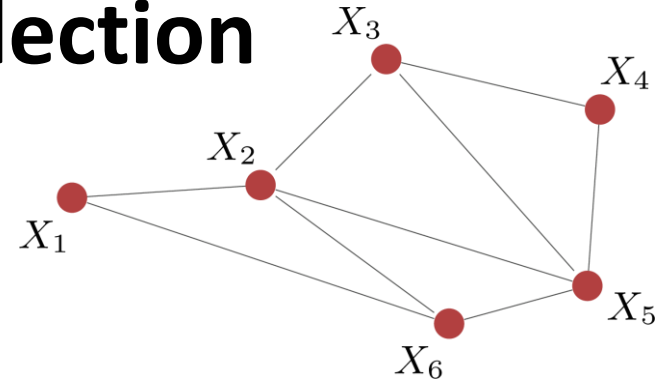
sample covariance matrix

$$\hat{\Sigma}_n := \frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)\top}$$

Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

$$\text{pdf: } f_X(x; \Theta^*) = \frac{\sqrt{\det \Theta^*}}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{1}{2}x^\top \Theta^* x\right)$$



*Maximum likelihood estimator of Θ^**

$$\hat{\Theta}_{\text{ML}} = \arg \max_{\Theta} \log \prod_{i=1}^n f(x^{(1)}, \dots, x^{(n)}; \Theta)$$

$$= \arg \min_{\Theta} -\log \det \Theta + \text{tr}\left(\Theta \hat{\Sigma}_n\right)$$

$$= \hat{\Sigma}_n^{-1}$$

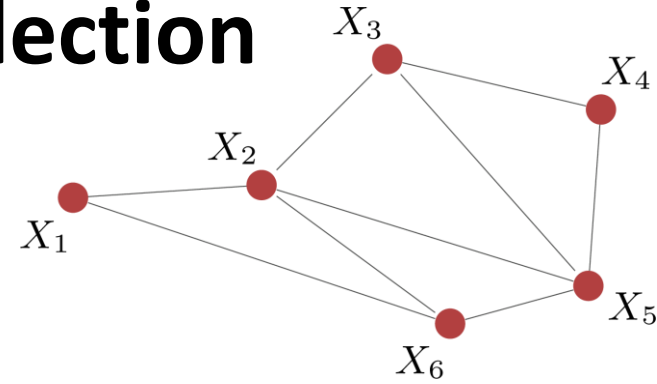
sample covariance matrix

$$\hat{\Sigma}_n := \frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)\top}$$

Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

$$\text{pdf: } f_X(x; \Theta^*) = \frac{\sqrt{\det \Theta^*}}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{1}{2}x^\top \Theta^* x\right)$$



*Maximum likelihood estimator of Θ^**

$$\hat{\Theta}_{\text{ML}} = \arg \max_{\Theta} \log \prod_{i=1}^n f(x^{(1)}, \dots, x^{(n)}; \Theta)$$

$$= \arg \min_{\Theta} -\log \det \Theta + \text{tr}\left(\Theta \hat{\Sigma}_n\right)$$

$$= \hat{\Sigma}_n^{-1}$$

sample covariance matrix

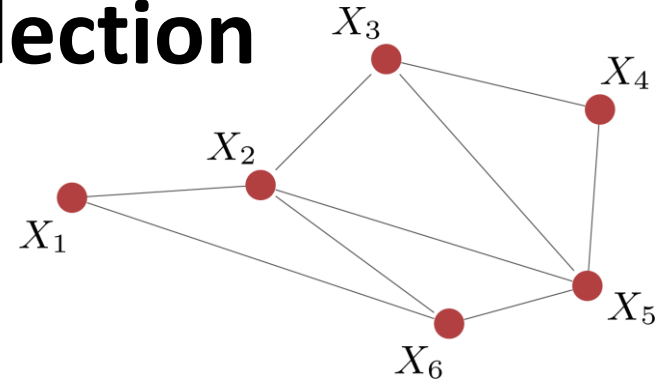
$$\hat{\Sigma}_n := \frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)\top}$$

(assuming it is invertible $\implies n > d$)

Gaussian Graphical Model Selection

$$X = (X_1, X_2, \dots, X_d) \sim \mathcal{N}(0_d, \Sigma^*)$$

$$\text{pdf: } f_X(x; \Theta^*) = \frac{\sqrt{\det \Theta^*}}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{1}{2}x^\top \Theta^* x\right)$$



*Maximum likelihood estimator of Θ^**

$$\hat{\Theta}_{\text{ML}} = \arg \max_{\Theta} \log \prod_{i=1}^n f(x^{(1)}, \dots, x^{(n)}; \Theta)$$

$$= \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n)$$

sample covariance matrix

$$= \hat{\Sigma}_n^{-1}$$

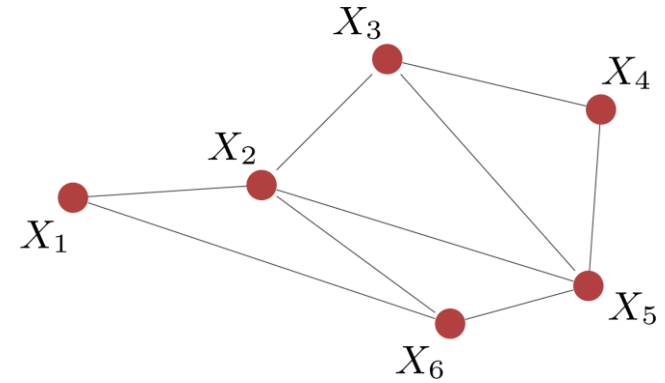
$$\hat{\Sigma}_n := \frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)\top}$$

(assuming it is invertible $\implies n > d$)

Graphical LASSO

Maximum likelihood estimator

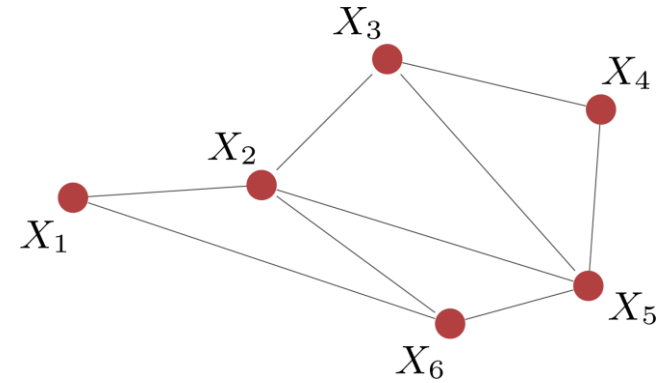
$$\hat{\Theta}_{\text{ML}}^* = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n)$$



Graphical LASSO

Maximum likelihood estimator

$$\hat{\Theta}_{\text{ML}}^* = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n)$$

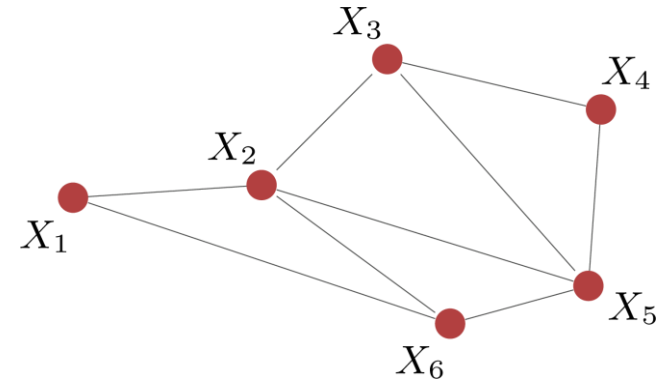


Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent

Graphical LASSO

Maximum likelihood estimator

$$\hat{\Theta}_{\text{ML}}^* = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n)$$



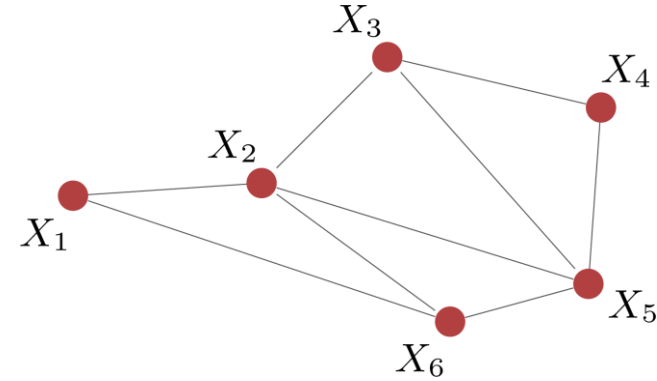
Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent

Graphical LASSO

Graphical LASSO

Maximum likelihood estimator

$$\hat{\Theta}_{\text{ML}}^* = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n)$$



Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent

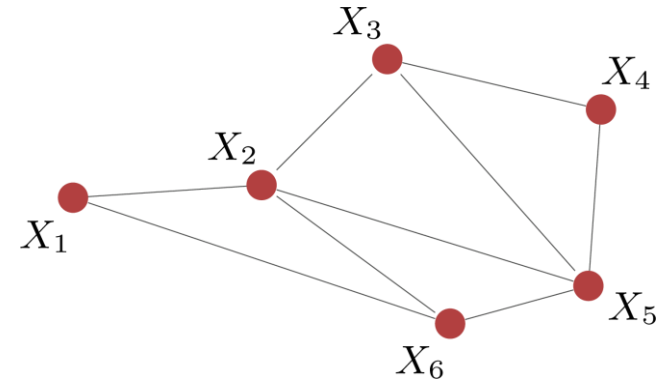
Graphical LASSO

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1,\text{off-d}}$$

Graphical LASSO

Maximum likelihood estimator

$$\hat{\Theta}_{\text{ML}}^* = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n)$$



Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent

Graphical LASSO

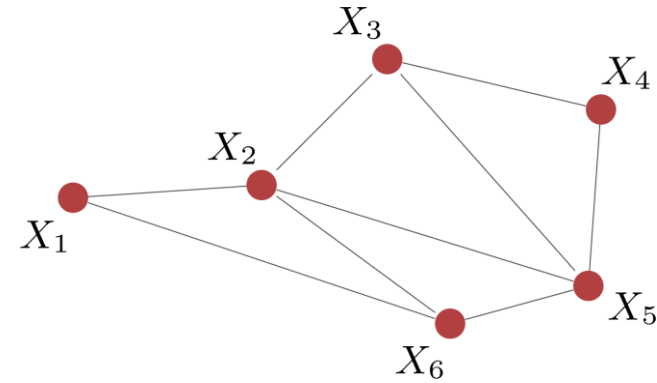
$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1, \text{off-d}}$$

$$\left| \sum_{i \neq j} |\Theta_{ij}| \right|$$

Graphical LASSO

Maximum likelihood estimator

$$\hat{\Theta}_{\text{ML}}^* = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n)$$



Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent

Graphical LASSO

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1,\text{off-d}}$$

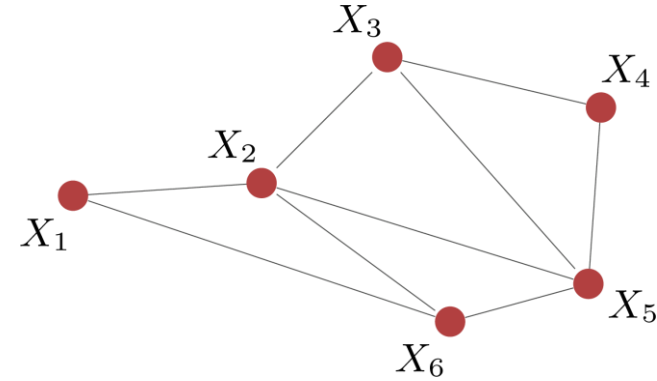
$$\sum_{i \neq j} |\Theta_{ij}|$$

applies L1-norm only to off-diagonal entries

Graphical LASSO

Maximum likelihood estimator

$$\hat{\Theta}_{\text{ML}}^* = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n)$$



Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent

Graphical LASSO

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1, \text{off-d}}$$

$$\sum_{i \neq j} |\Theta_{ij}|$$

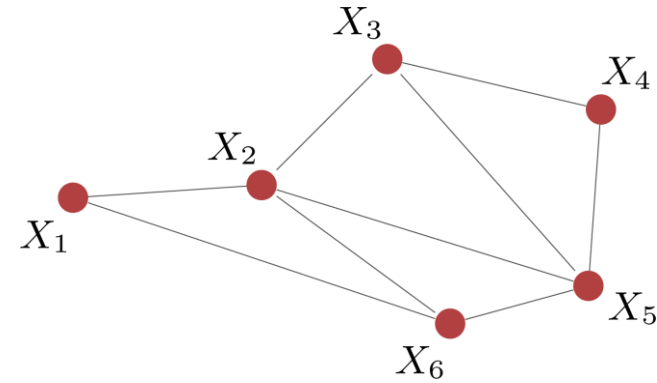
applies L1-norm only to off-diagonal entries

$$\Theta^* = \begin{bmatrix} * & * & 0 & 0 & 0 & * \\ * & * & * & 0 & * & * \\ 0 & * & * & * & * & 0 \\ 0 & 0 & * & * & * & 0 \\ 0 & * & * & * & * & * \\ * & * & 0 & 0 & * & * \end{bmatrix}$$

Graphical LASSO

Maximum likelihood estimator

$$\hat{\Theta}_{\text{ML}}^* = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n)$$



Assumption: most pairs of coordinates (X_i, X_j) are conditionally independent

Graphical LASSO

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1, \text{off-d}}$$

$$\sum_{i \neq j} |\Theta_{ij}|$$

applies L1-norm only to off-diagonal entries

$$\Theta^* = \begin{bmatrix} * & * & 0 & 0 & 0 & * \\ * & * & * & 0 & * & * \\ 0 & * & * & * & * & 0 \\ 0 & 0 & * & * & * & 0 \\ 0 & * & * & * & * & * \\ * & * & 0 & 0 & * & * \end{bmatrix}$$

sensible estimators even for non-Gaussian RVs

Example: Dow Jones

Example: Dow Jones

Price of stock of *6 companies* at beginning of each week of 2011 (Jan-Jun)

Example: Dow Jones

Price of stock of *6 companies* at beginning of each week of 2011 (Jan-Jun)

Week	Alcoa	American Express	Boeing	Bank of America	Caterpillar	Cisco Systems
1	14.67	43.30	66.15	10.59	100.25	14.94
2	15.29	43.73	69.26	10.89	101.30	15.14
3	15.82	43.86	69.42	11.18	102.59	16.04
4	15.87	43.86	70.29	11.47	102.72	16.41
5	15.92	43.96	70.86	11.87	103.42	16.59
6	15.95	44.13	71.17	11.89	103.56	16.82
7	15.96	44.20	71.43	12.28	104.86	16.88
8	16.18	44.75	71.52	12.32	105.58	16.93
9	16.19	44.94	71.60	12.36	105.87	17.01
⋮	⋮	⋮	⋮	⋮	⋮	⋮
24	17.42	50.74	79.31	14.77	96.93	21.22
25	18.06	51.39	80.35	15.08	99.62	22.11

Example: Dow Jones

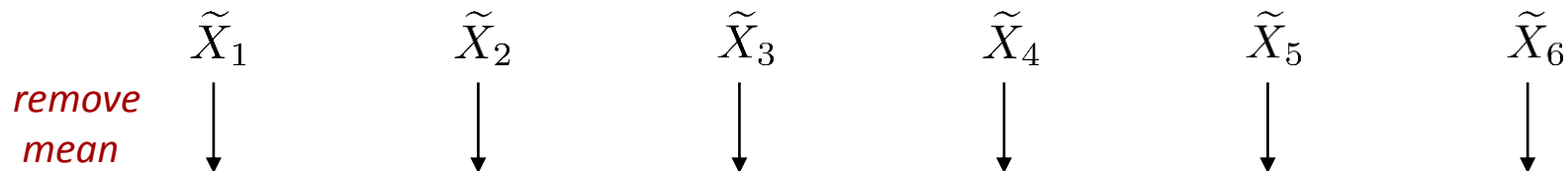
Price of stock of *6 companies* at beginning of each week of 2011 (Jan-Jun)

Week	Alcoa	American Express	Boeing	Bank of America	Caterpillar	Cisco Systems
1	14.67	43.30	66.15	10.59	100.25	14.94
2	15.29	43.73	69.26	10.89	101.30	15.14
3	15.82	43.86	69.42	11.18	102.59	16.04
4	15.87	43.86	70.29	11.47	102.72	16.41
5	15.92	43.96	70.86	11.87	103.42	16.59
6	15.95	44.13	71.17	11.89	103.56	16.82
7	15.96	44.20	71.43	12.28	104.86	16.88
8	16.18	44.75	71.52	12.32	105.58	16.93
9	16.19	44.94	71.60	12.36	105.87	17.01
⋮	⋮	⋮	⋮	⋮	⋮	⋮
24	17.42	50.74	79.31	14.77	96.93	21.22
25	18.06	51.39	80.35	15.08	99.62	22.11
	\tilde{X}_1	\tilde{X}_2	\tilde{X}_3	\tilde{X}_4	\tilde{X}_5	\tilde{X}_6

Example: Dow Jones

Price of stock of *6 companies* at beginning of each week of 2011 (Jan-Jun)

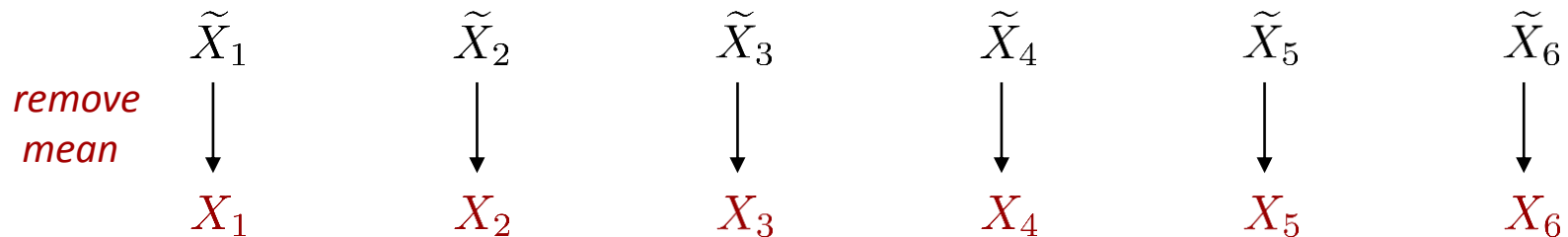
Week	Alcoa	American Express	Boeing	Bank of America	Caterpillar	Cisco Systems
1	14.67	43.30	66.15	10.59	100.25	14.94
2	15.29	43.73	69.26	10.89	101.30	15.14
3	15.82	43.86	69.42	11.18	102.59	16.04
4	15.87	43.86	70.29	11.47	102.72	16.41
5	15.92	43.96	70.86	11.87	103.42	16.59
6	15.95	44.13	71.17	11.89	103.56	16.82
7	15.96	44.20	71.43	12.28	104.86	16.88
8	16.18	44.75	71.52	12.32	105.58	16.93
9	16.19	44.94	71.60	12.36	105.87	17.01
⋮	⋮	⋮	⋮	⋮	⋮	⋮
24	17.42	50.74	79.31	14.77	96.93	21.22
25	18.06	51.39	80.35	15.08	99.62	22.11



Example: Dow Jones

Price of stock of *6 companies* at beginning of each week of 2011 (Jan-Jun)

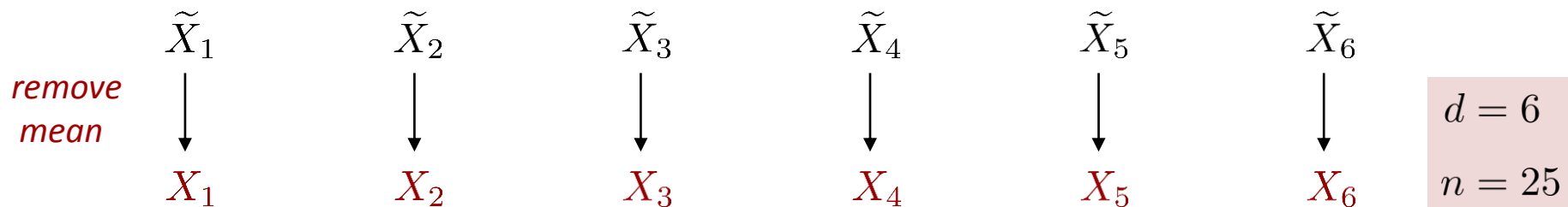
Week	Alcoa	American Express	Boeing	Bank of America	Caterpillar	Cisco Systems
1	14.67	43.30	66.15	10.59	100.25	14.94
2	15.29	43.73	69.26	10.89	101.30	15.14
3	15.82	43.86	69.42	11.18	102.59	16.04
4	15.87	43.86	70.29	11.47	102.72	16.41
5	15.92	43.96	70.86	11.87	103.42	16.59
6	15.95	44.13	71.17	11.89	103.56	16.82
7	15.96	44.20	71.43	12.28	104.86	16.88
8	16.18	44.75	71.52	12.32	105.58	16.93
9	16.19	44.94	71.60	12.36	105.87	17.01
⋮	⋮	⋮	⋮	⋮	⋮	⋮
24	17.42	50.74	79.31	14.77	96.93	21.22
25	18.06	51.39	80.35	15.08	99.62	22.11



Example: Dow Jones

Price of stock of *6 companies* at beginning of each week of 2011 (Jan-Jun)

Week	Alcoa	American Express	Boeing	Bank of America	Caterpillar	Cisco Systems
1	14.67	43.30	66.15	10.59	100.25	14.94
2	15.29	43.73	69.26	10.89	101.30	15.14
3	15.82	43.86	69.42	11.18	102.59	16.04
4	15.87	43.86	70.29	11.47	102.72	16.41
5	15.92	43.96	70.86	11.87	103.42	16.59
6	15.95	44.13	71.17	11.89	103.56	16.82
7	15.96	44.20	71.43	12.28	104.86	16.88
8	16.18	44.75	71.52	12.32	105.58	16.93
9	16.19	44.94	71.60	12.36	105.87	17.01
⋮	⋮	⋮	⋮	⋮	⋮	⋮
24	17.42	50.74	79.31	14.77	96.93	21.22
25	18.06	51.39	80.35	15.08	99.62	22.11



Example: Dow Jones

Example: Dow Jones

$$\hat{\Theta}_{\text{ML}} = \hat{\Sigma}_n^{-1} = \begin{bmatrix} 60.82 & 4.85 & -6.21 & -21.34 & -0.07 & -4.81 \\ 4.85 & 7.34 & -1.22 & -5.50 & 0.37 & -4.78 \\ -6.21 & -1.22 & 3.03 & 2.08 & -0.08 & -2.69 \\ -21.34 & -5.50 & 2.08 & 14.31 & -0.42 & 1.69 \\ -0.07 & 0.37 & -0.08 & -0.42 & 0.06 & 0.07 \\ -4.81 & -4.78 & -2.69 & 1.69 & 0.07 & 11.38 \end{bmatrix}$$

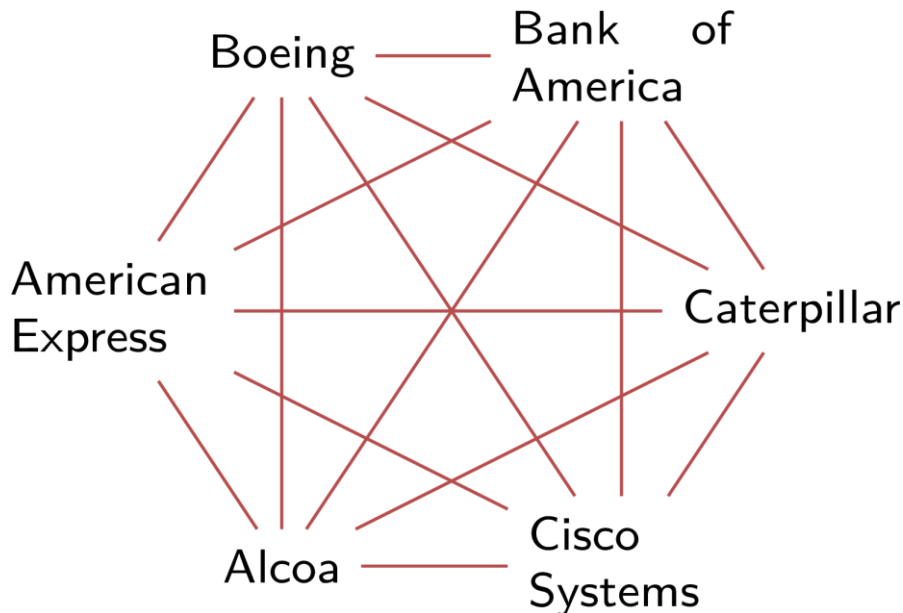
Example: Dow Jones

$$\hat{\Theta}_{\text{ML}} = \hat{\Sigma}_n^{-1} = \begin{bmatrix} 60.82 & 4.85 & -6.21 & -21.34 & -0.07 & -4.81 \\ 4.85 & 7.34 & -1.22 & -5.50 & 0.37 & -4.78 \\ -6.21 & -1.22 & 3.03 & 2.08 & -0.08 & -2.69 \\ -21.34 & -5.50 & 2.08 & 14.31 & -0.42 & 1.69 \\ -0.07 & 0.37 & -0.08 & -0.42 & 0.06 & 0.07 \\ -4.81 & -4.78 & -2.69 & 1.69 & 0.07 & 11.38 \end{bmatrix}$$

Boeing	Bank of America
American Express	Caterpillar
Alcoa	Cisco Systems

Example: Dow Jones

$$\hat{\Theta}_{\text{ML}} = \hat{\Sigma}_n^{-1} = \begin{bmatrix} 60.82 & 4.85 & -6.21 & -21.34 & -0.07 & -4.81 \\ 4.85 & 7.34 & -1.22 & -5.50 & 0.37 & -4.78 \\ -6.21 & -1.22 & 3.03 & 2.08 & -0.08 & -2.69 \\ -21.34 & -5.50 & 2.08 & 14.31 & -0.42 & 1.69 \\ -0.07 & 0.37 & -0.08 & -0.42 & 0.06 & 0.07 \\ -4.81 & -4.78 & -2.69 & 1.69 & 0.07 & 11.38 \end{bmatrix}$$



Example: Dow Jones

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1, \text{off-d}}$$

graphical LASSO

Example: Dow Jones

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1, \text{off-d}} \quad \textit{graphical LASSO}$$

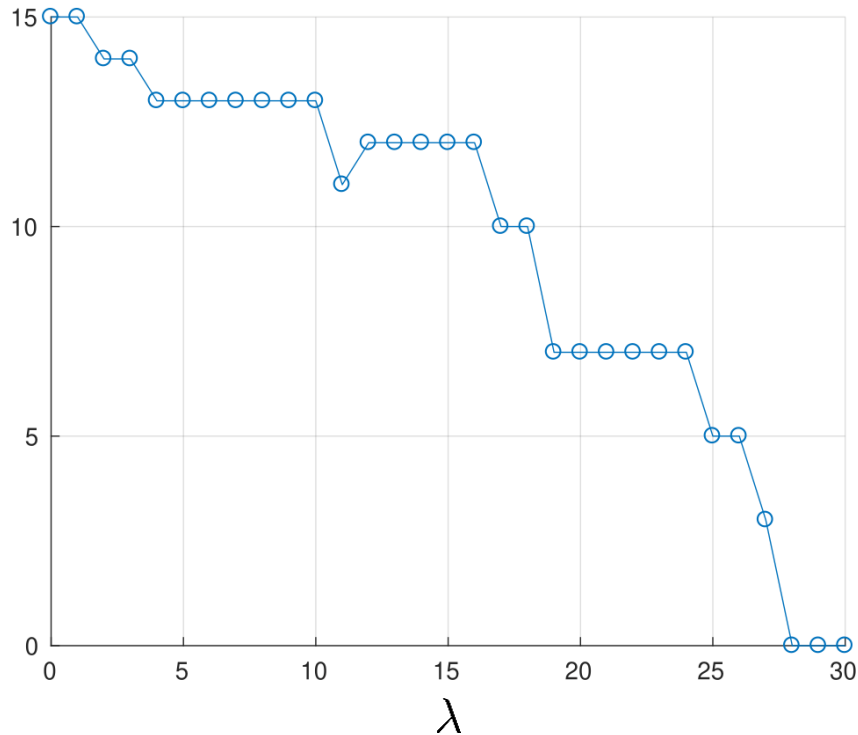
$|(\hat{\Theta}_{\text{GL}})_{ij}| \leq 10^{-3} \implies$ we assume no correlation, i.e., *no edge* (i, j)

Example: Dow Jones

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1, \text{off-d}} \quad \text{graphical LASSO}$$

$|(\hat{\Theta}_{\text{GL}})_{ij}| \leq 10^{-3} \implies$ we assume no correlation, i.e., *no edge* (i, j)

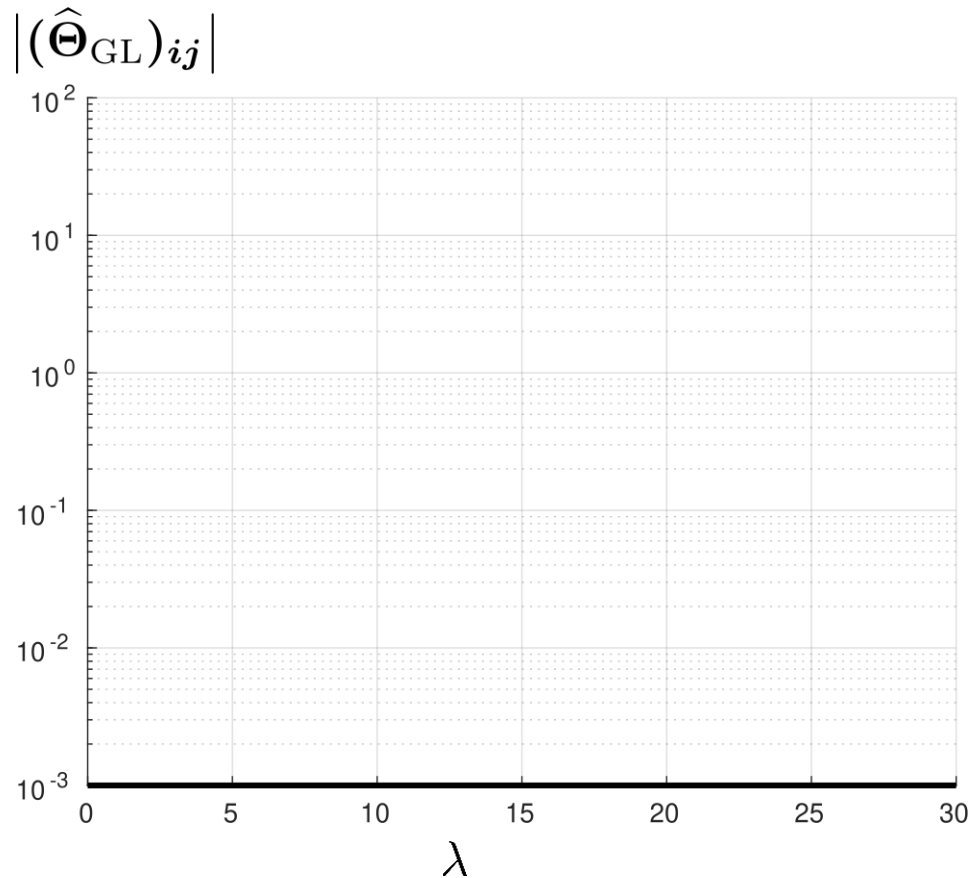
Number of edges



Example: Dow Jones

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1, \text{off-d}} \quad \text{graphical LASSO}$$

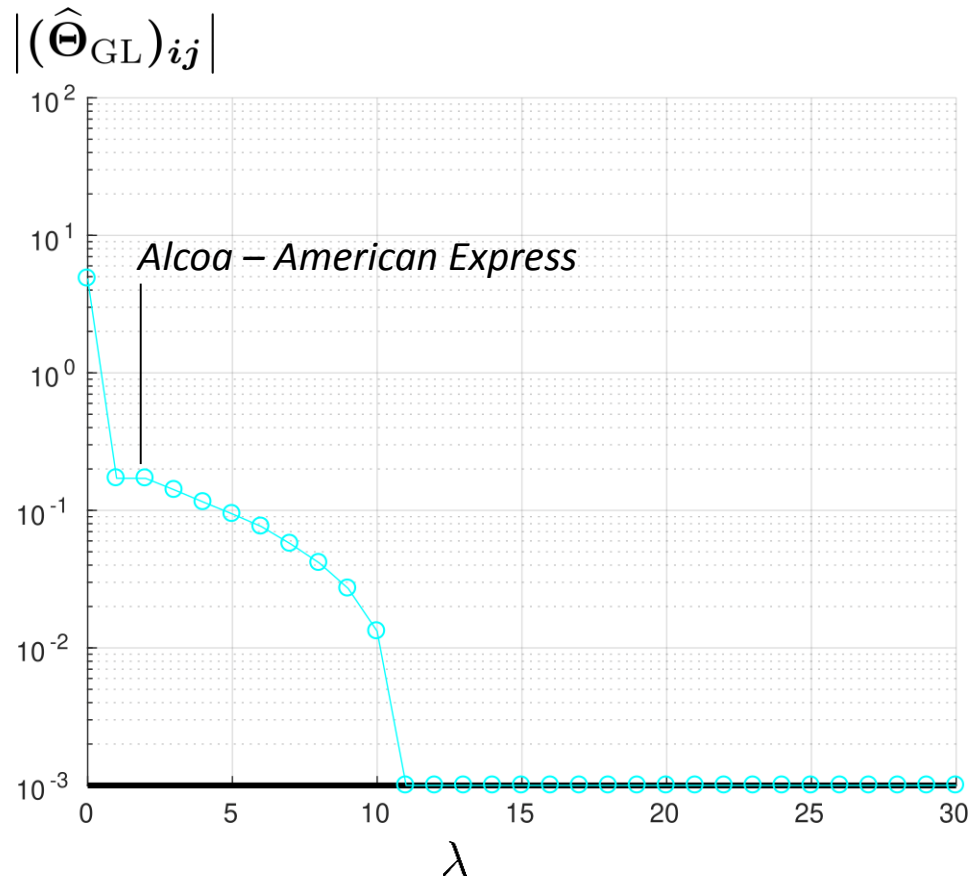
$|(\hat{\Theta}_{\text{GL}})_{ij}| \leq 10^{-3} \implies$ we assume no correlation, i.e., *no edge* (i, j)



Example: Dow Jones

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1,\text{off-d}} \quad \text{graphical LASSO}$$

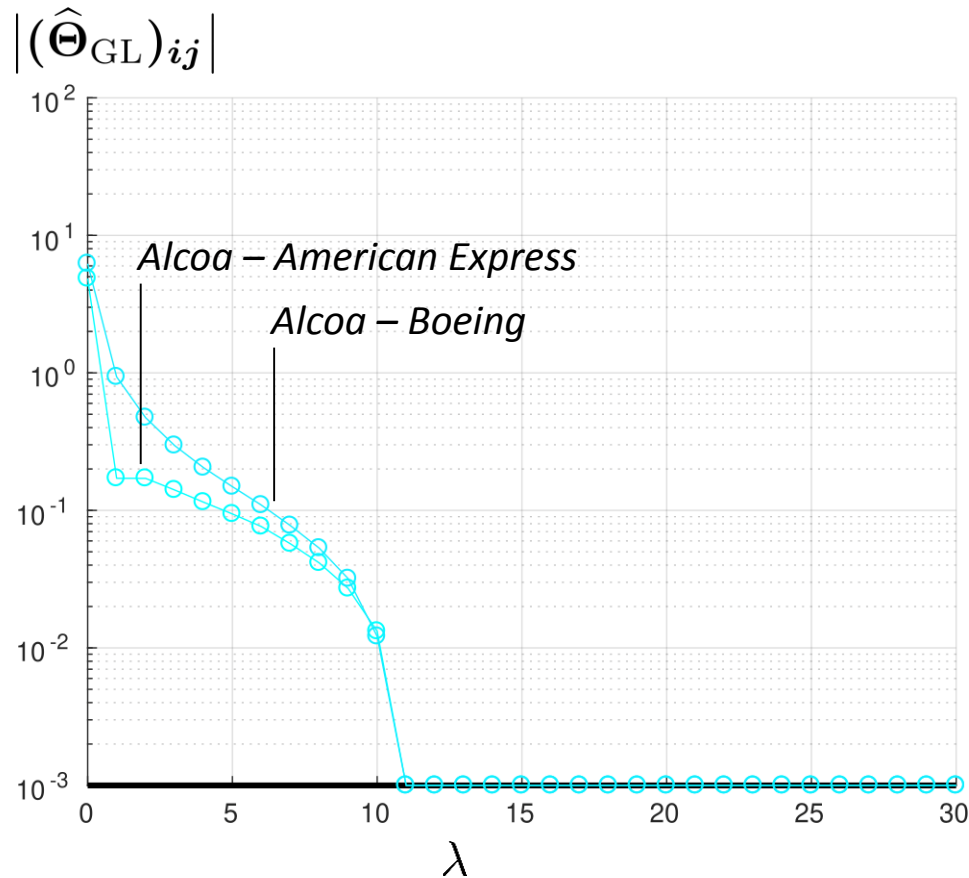
$|(\hat{\Theta}_{\text{GL}})_{ij}| \leq 10^{-3} \implies$ we assume no correlation, i.e., *no edge* (i, j)



Example: Dow Jones

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1,\text{off-d}} \quad \text{graphical LASSO}$$

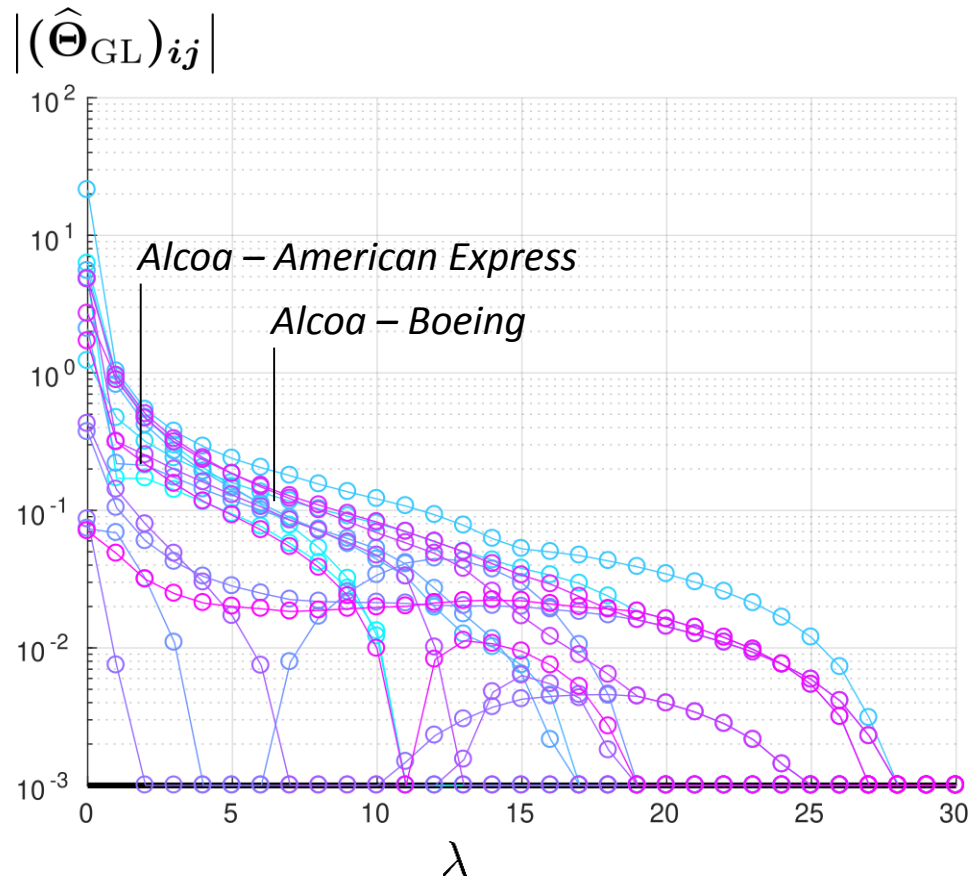
$|(\hat{\Theta}_{\text{GL}})_{ij}| \leq 10^{-3} \implies$ we assume no correlation, i.e., *no edge* (i, j)



Example: Dow Jones

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1,\text{off-d}} \quad \text{graphical LASSO}$$

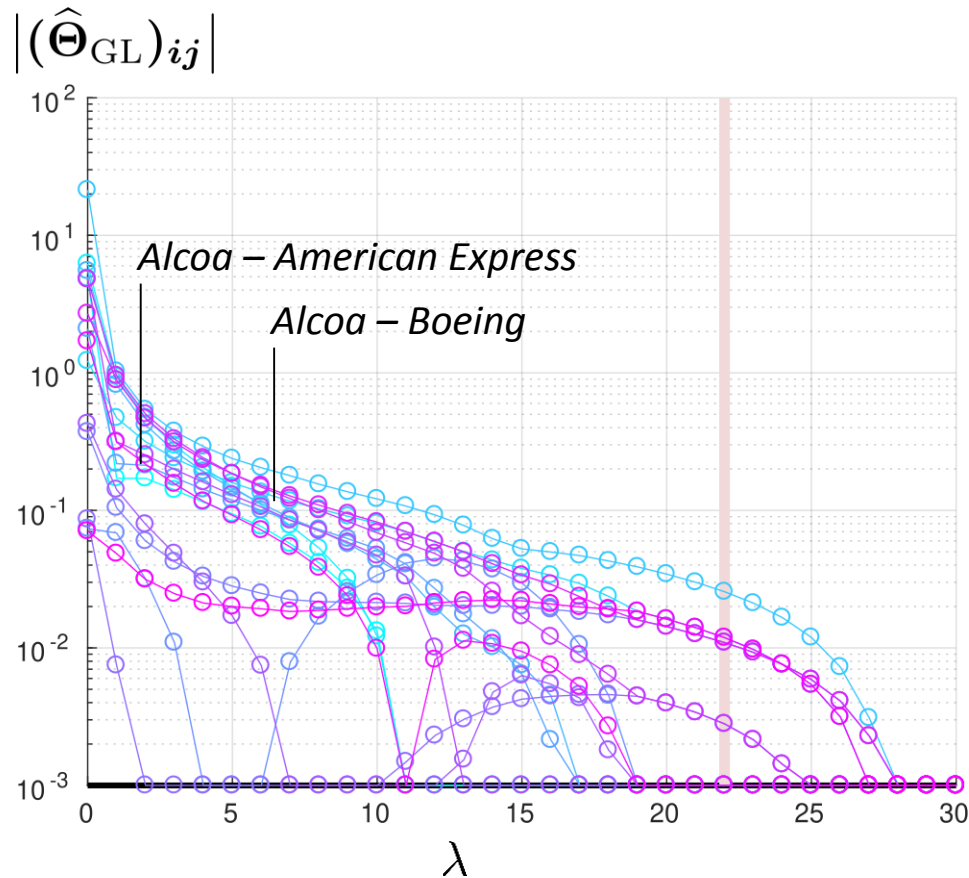
$|(\hat{\Theta}_{\text{GL}})_{ij}| \leq 10^{-3} \implies$ we assume no correlation, i.e., *no edge* (i, j)



Example: Dow Jones

$$\hat{\Theta}_{\text{GL}} = \arg \min_{\Theta} -\log \det \Theta + \text{tr}(\Theta \hat{\Sigma}_n) + \lambda \|\Theta\|_{1,\text{off-d}} \quad \text{graphical LASSO}$$

$|(\hat{\Theta}_{\text{GL}})_{ij}| \leq 10^{-3} \implies$ we assume no correlation, i.e., *no edge* (i, j)



Example: Dow Jones

$$\hat{\Theta}_{\text{GL}} = \begin{bmatrix} 1.814 & 0 & 0 & -0.030 & 0 & 0 \\ 0 & 0.186 & -0.014 & 0 & 0.013 & -0.003 \\ 0 & -0.014 & 0.096 & 0 & 0.003 & -0.013 \\ -0.030 & 0 & 0 & 0.601 & 0 & 0 \\ 0 & 0.013 & 0.003 & 0 & 0.026 & 0.014 \\ 0 & -0.003 & -0.013 & 0 & 0.014 & 0.282 \end{bmatrix}$$

Example: Dow Jones

$$\hat{\Theta}_{GL} = \begin{bmatrix} 1.814 & 0 & 0 & -0.030 & 0 & 0 \\ 0 & 0.186 & -0.014 & 0 & 0.013 & -0.003 \\ 0 & -0.014 & 0.096 & 0 & 0.003 & -0.013 \\ -0.030 & 0 & 0 & 0.601 & 0 & 0 \\ 0 & 0.013 & 0.003 & 0 & 0.026 & 0.014 \\ 0 & -0.003 & -0.013 & 0 & 0.014 & 0.282 \end{bmatrix}$$

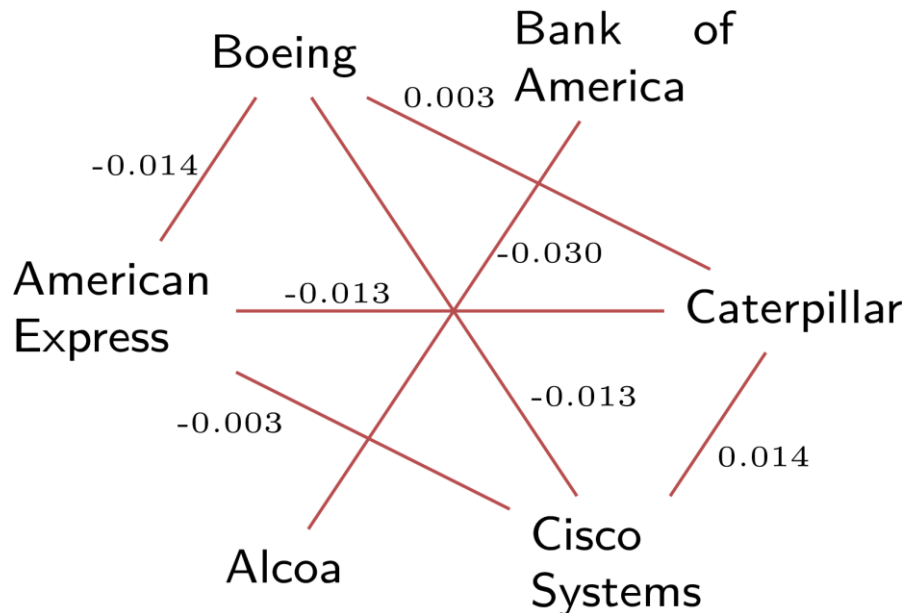
Boeing Bank of
America

American Caterpillar
Express

Alcoa Cisco
Systems

Example: Dow Jones

$$\hat{\Theta}_{GL} = \begin{bmatrix} 1.814 & 0 & 0 & -0.030 & 0 & 0 \\ 0 & 0.186 & -0.014 & 0 & 0.013 & -0.003 \\ 0 & -0.014 & 0.096 & 0 & 0.003 & -0.013 \\ -0.030 & 0 & 0 & 0.601 & 0 & 0 \\ 0 & 0.013 & 0.003 & 0 & 0.026 & 0.014 \\ 0 & -0.003 & -0.013 & 0 & 0.014 & 0.282 \end{bmatrix}$$



Outline

- *Motivation: Hypothesis Testing in High-Dimensions*
- *Introduction to LASSO and other sparsity problems*
- *Gaussian graphical model selection*
- Matrix completion

Matrix Completion

Matrix Completion

Suppose someone gave you \$1M for *completing* a table like this...

Matrix Completion

Suppose someone gave you \$1M for *completing* a table like this...

		movies							
users	1	★	0.1	★	★	0.9	★	0.9	★
	2	0.5	0.4	★	★	★	0.6	★	0.6
	3	★	★	★	0.7	★	★	★	0.2
	4	0.4	0.2	0.9	★	★	★	★	★
	5	★	★	0.7	0.7	★	★	0.3	★

Matrix Completion

Suppose someone gave you \$1M for *completing* a table like this...

		movies							
users	1	★	0.1	★	★	0.9	★	0.9	★
	2	0.5	0.4	★	★	★	0.6	★	0.6
	3	★	★	★	0.7	★	★	★	0.2
	4	0.4	0.2	0.9	★	★	★	★	★
	5	★	★	0.7	0.7	★	★	0.3	★

— rating user i gives to movie j

Matrix Completion

Suppose someone gave you \$1M for *completing* a table like this...

		movies							
users	1	★	0.1	★	★	0.9	★	0.9	★
	2	0.5	0.4	★	★	★	0.6	★	0.6
	3	★	★	★	0.7	★	★	★	0.2
	4	0.4	0.2	0.9	★	★	★	★	★
	5	★	★	0.7	0.7	★	★	0.3	★

— rating user i gives to movie j

Key insight:

Matrix Completion

Suppose someone gave you \$1M for *completing* a table like this...

		movies							
users	1	★	0.1	★	★	0.9	★	0.9	★
	2	0.5	0.4	★	★	★	0.6	★	0.6
	3	★	★	★	0.7	★	★	★	0.2
	4	0.4	0.2	0.9	★	★	★	★	★
	5	★	★	0.7	0.7	★	★	0.3	★

— rating user i gives to movie j

Key insight: only a few factors may explain users' tastes (genre, actors, ads, ...)

Matrix Completion

Suppose someone gave you \$1M for *completing* a table like this...

		movies							
users		★	0.1	★	★	0.9	★	0.9	★
		0.5	0.4	★	★	★	0.6	★	0.6
		★	★	★	0.7	★	★	★	0.2
		0.4	0.2	0.9	★	★	★	★	★
		★	★	0.7	0.7	★	★	0.3	★

— rating user i gives to movie j

Key insight: only a few factors may explain users' tastes (genre, actors, ads, ...)

ratings

X^*

Matrix Completion

Suppose someone gave you \$1M for *completing* a table like this...

		movies						
users	★	0.1	★	★	0.9	★	0.9	★
	0.5	0.4	★	★	★	0.6	★	0.6
	★	★	★	0.7	★	★	★	0.2
	0.4	0.2	0.9	★	★	★	★	★
	★	★	0.7	0.7	★	★	0.3	★

rating user i gives to movie j

Key insight: only a few factors may explain users' tastes (genre, actors, ads, ...)

ratings

$$\begin{array}{|c|} \hline X^* \\ \hline \end{array} = \begin{array}{|c|} \hline \tilde{U} \\ \hline \end{array} \begin{array}{|c|} \hline \tilde{V} \\ \hline \end{array}$$

Matrix Completion

Suppose someone gave you \$1M for *completing* a table like this...

		movies						
users	★	0.1	★	★	0.9	★	0.9	★
	0.5	0.4	★	★	★	0.6	★	0.6
	★	★	★	0.7	★	★	★	0.2
	0.4	0.2	0.9	★	★	★	★	★
	★	★	0.7	0.7	★	★	0.3	★

rating user i gives to movie j

Key insight: only a few factors may explain users' tastes (genre, actors, ads, ...)

ratings

$$\begin{array}{c} \text{ratings} \\ X^* \end{array} = \begin{array}{c} \tilde{U} \\ \text{influence of factors on users' tastes} \end{array} \begin{array}{c} \tilde{V} \end{array}$$

influence of factors on users' tastes

Matrix Completion

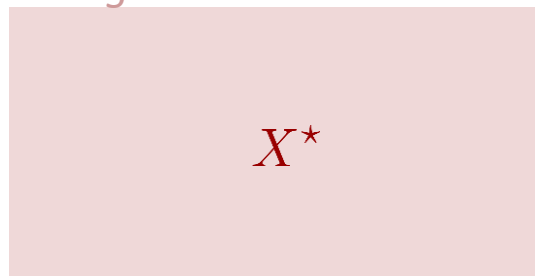
Suppose someone gave you \$1M for *completing* a table like this...

		movies						
users	★	0.1	★	★	0.9	★	0.9	★
	0.5	0.4	★	★	★	0.6	★	0.6
	★	★	★	0.7	★	★	★	0.2
	0.4	0.2	0.9	★	★	★	★	★
	★	★	0.7	0.7	★	★	0.3	★

rating user i gives to movie j

Key insight: only a few factors may explain users' tastes (genre, actors, ads, ...)

ratings



=



influence of factors on users' tastes

characterization of movies



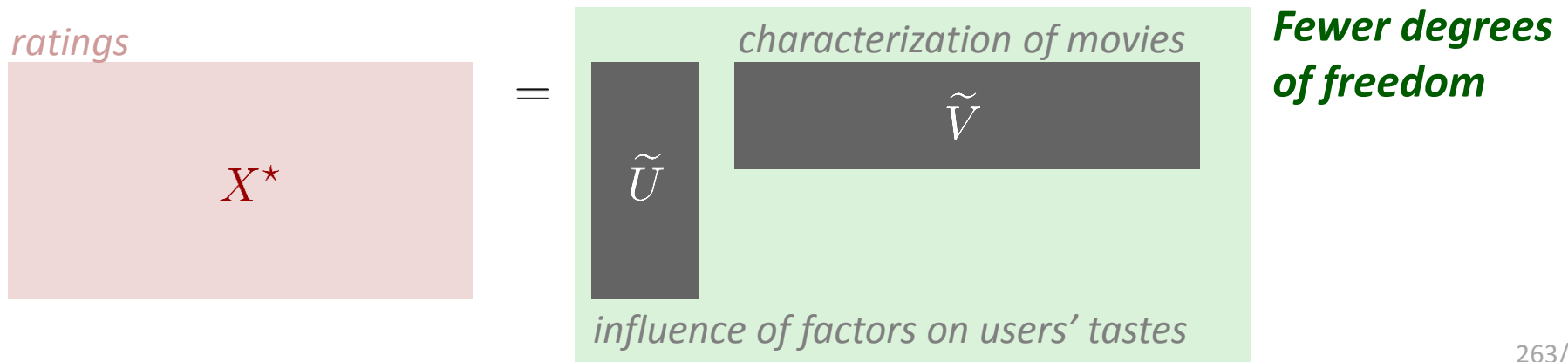
Matrix Completion

Suppose someone gave you \$1M for *completing* a table like this...

		movies						
users	★	0.1	★	★	0.9	★	0.9	★
	0.5	0.4	★	★	★	0.6	★	0.6
	★	★	★	0.7	★	★	★	0.2
	0.4	0.2	0.9	★	★	★	★	★
	★	★	0.7	0.7	★	★	0.3	★

— rating user i gives to movie j

Key insight: only a few factors may explain users' tastes (genre, actors, ads, ...)



Singular value decomposition:

Singular value decomposition: any real $m \times n$ matrix can be decomposed as

Singular value decomposition: any real $m \times n$ matrix can be decomposed as

$$A = U\Sigma V^{\top}$$

Singular value decomposition: any real $m \times n$ matrix can be decomposed as

$$A = U\Sigma V^{\top} = \underbrace{\begin{bmatrix} | & & | \\ u_1 & \cdots & u_k \\ | & & | \end{bmatrix}}_{m \times k} \underbrace{\begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_k \end{bmatrix}}_{k \times k} \underbrace{\begin{bmatrix} - & v_1^{\top} & - \\ & \vdots & \\ - & v_k^{\top} & - \end{bmatrix}}_{k \times n}$$

Singular value decomposition: any real $m \times n$ matrix can be decomposed as

$$A = U\Sigma V^{\top} = \underbrace{\begin{bmatrix} | & & | \\ u_1 & \cdots & u_k \\ | & & | \end{bmatrix}}_{\substack{m \times k \\ \text{orthogonal}}} \underbrace{\begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_k \end{bmatrix}}_{k \times k} \underbrace{\begin{bmatrix} - & v_1^{\top} & - \\ & \vdots & \\ - & v_k^{\top} & - \end{bmatrix}}_{\substack{k \times n \\ \text{orthogonal}}}$$

Singular value decomposition: any real $m \times n$ matrix can be decomposed as

$$A = U\Sigma V^{\top} = \underbrace{\begin{bmatrix} | & & | \\ u_1 & \cdots & u_k \\ | & & | \end{bmatrix}}_{\substack{m \times k \\ \text{orthogonal}}} \underbrace{\begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_k \end{bmatrix}}_{\substack{k \times k \\ \text{rank}}} \underbrace{\begin{bmatrix} - & v_1^{\top} & - \\ & \vdots & \\ - & v_k^{\top} & - \end{bmatrix}}_{\substack{k \times n \\ \text{orthogonal}}}$$

Singular value decomposition: any real $m \times n$ matrix can be decomposed as

$$\begin{aligned}
 A = U\Sigma V^\top &= \underbrace{\begin{bmatrix} | & & | \\ u_1 & \cdots & u_k \\ | & & | \end{bmatrix}}_{\substack{m \times k \\ \text{orthogonal}}} \underbrace{\begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_k \end{bmatrix}}_{\substack{k \times k \\ \text{rank}}} \underbrace{\begin{bmatrix} - & v_1^\top & - \\ & \vdots & \\ - & v_k^\top & - \end{bmatrix}}_{\substack{k \times n \\ \text{orthogonal}}} \\
 &= \underbrace{\begin{bmatrix} | & & | & | & \\ u_1 & \cdots & u_k & u'_{k+1} & \cdots & u'_m \\ | & & | & | & \end{bmatrix}}_{m \times m} \underbrace{\left[\begin{array}{cccc|c} \sigma_1 & 0 & \cdots & 0 & \\ 0 & \sigma_2 & \cdots & 0 & \\ \vdots & \vdots & \ddots & \vdots & \\ 0 & 0 & \cdots & \sigma_k & \\ \hline & & & & 0_{k \times (n-k)} \end{array} \right]}_{\substack{m \times n}} \underbrace{\begin{bmatrix} - & v_1^\top & - \\ & \vdots & \\ - & v_k^\top & - \\ & \vdots & \\ - & v_k'^\top & - \end{bmatrix}}_{n \times n}
 \end{aligned}$$

Singular value decomposition: any real $m \times n$ matrix can be decomposed as

$$A = U\Sigma V^{\top} = \underbrace{\begin{bmatrix} | & & | \\ u_1 & \cdots & u_k \\ | & & | \end{bmatrix}}_{\substack{m \times k \\ \text{orthogonal}}} \underbrace{\begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_k \end{bmatrix}}_{\substack{k \times k \\ \text{rank}}} \underbrace{\begin{bmatrix} - & v_1^{\top} & - \\ & \vdots & \\ - & v_k^{\top} & - \end{bmatrix}}_{\substack{k \times n \\ \text{orthogonal}}}$$

$$= \underbrace{\begin{bmatrix} | & & | & | & \\ u_1 & \cdots & u_k & u'_{k+1} & \cdots & u'_m \\ | & & | & | & \\ \hline \end{bmatrix}}_{m \times m} \underbrace{\left[\begin{array}{cccc|c} \sigma_1 & 0 & \cdots & 0 & \\ 0 & \sigma_2 & \cdots & 0 & \\ \vdots & \vdots & \ddots & \vdots & \\ 0 & 0 & \cdots & \sigma_k & \\ \hline \end{array} \right]}_{m \times n} \underbrace{\begin{bmatrix} - & v_1^{\top} & - \\ & \vdots & \\ - & v_k^{\top} & - \\ & v_1'^{\top} & - \\ & \vdots & \\ - & v_k'^{\top} & - \end{bmatrix}}_{n \times n}$$

Singular value decomposition: any real $m \times n$ matrix can be decomposed as

$$A = U\Sigma V^{\top} = \underbrace{\begin{bmatrix} | & & | \\ u_1 & \cdots & u_k \\ | & & | \end{bmatrix}}_{\substack{m \times k \\ \text{orthogonal}}} \underbrace{\begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_k \end{bmatrix}}_{\substack{k \times k \\ \text{rank}}} \underbrace{\begin{bmatrix} - & v_1^{\top} & - \\ & \vdots & \\ - & v_k^{\top} & - \end{bmatrix}}_{\substack{k \times n \\ \text{orthogonal}}}$$

$$= \underbrace{\begin{bmatrix} | & \cdots & | & | & \cdots & | \\ u_1 & \cdots & u_k & u'_{k+1} & \cdots & u'_m \\ | & \cdots & | & | & \cdots & | \end{bmatrix}}_{m \times m} \underbrace{\left[\begin{array}{cccc|ccc} \sigma_1 & 0 & \cdots & 0 & & & \\ 0 & \sigma_2 & \cdots & 0 & & & \\ \vdots & \vdots & \ddots & \vdots & & & \\ 0 & 0 & \cdots & \sigma_k & & & \\ \hline & & & & 0_{k \times (n-k)} & & \\ & & & & 0_{(m-k) \times k} & 0_{(m-k) \times (n-k)} & \end{array} \right]}_{m \times n} \underbrace{\begin{bmatrix} - & v_1^{\top} & - \\ & \vdots & \\ - & v_k^{\top} & - \\ & v_1'^{\top} & - \\ & \vdots & \\ - & v_k'^{\top} & - \end{bmatrix}}_{n \times n}$$

sparse vector if matrix is low-rank

Singular value decomposition: any real $m \times n$ matrix can be decomposed as

$$A = U\Sigma V^{\top} = \underbrace{\begin{bmatrix} | & & | \\ u_1 & \cdots & u_k \\ | & & | \end{bmatrix}}_{\substack{m \times k \\ \text{orthogonal}}} \underbrace{\begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_k \end{bmatrix}}_{\substack{k \times k \\ \text{rank}}} \underbrace{\begin{bmatrix} - & v_1^{\top} & - \\ & \vdots & \\ - & v_k^{\top} & - \end{bmatrix}}_{\substack{k \times n \\ \text{orthogonal}}}$$

$$= \underbrace{\begin{bmatrix} | & & | & | & \\ u_1 & \cdots & u_k & u'_{k+1} & \cdots & u'_m \\ | & & | & | & \end{bmatrix}}_{m \times m} \underbrace{\left[\begin{array}{cccc|c} \sigma_1 & 0 & \cdots & 0 & \\ 0 & \sigma_2 & \cdots & 0 & \\ \vdots & \vdots & \ddots & \vdots & \\ 0 & 0 & \cdots & \sigma_k & \\ \hline & & & & 0_{k \times (n-k)} \\ \hline & & & & \\ \hline & & & & 0_{(m-k) \times (n-k)} \end{array} \right]}_{\substack{m \times n}} \underbrace{\begin{bmatrix} - & v_1^{\top} & - \\ & \vdots & \\ - & v_k^{\top} & - \\ & v_1'^{\top} & \\ & \vdots & \\ - & v_k'^{\top} & - \end{bmatrix}}_{n \times n}$$

$$= \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^{\top}}_{k \times n}$$

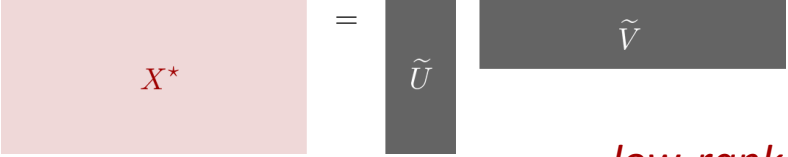
sparse vector if matrix is low-rank

$$X^{\star} = \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^{\top}}_{k \times n}$$

$$X^{\star} = \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^{\top}}_{k \times n}$$

$$\boxed{X^{\star}} = \boxed{\tilde{U}} \boxed{\tilde{V}}$$

$$X^{\star} = \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^{\top}}_{k \times n}$$



The diagram illustrates the low-rank decomposition of matrix X^{\star} . On the left is a light red rectangular block labeled X^{\star} . To its right is an equals sign. Further right is a dark gray vertical rectangular block labeled \tilde{U} . To the right of \tilde{U} is another dark gray horizontal rectangular block labeled \tilde{V} . Below the \tilde{V} block, the text *low-rank* is written in red.

$$X^{\star} = \tilde{U} \tilde{V}$$

low-rank

$$X^* = \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^{\top}}_{k \times n}$$

$$\boxed{X^*} = \boxed{\tilde{U}} \boxed{\tilde{V}}$$

low-rank

Our problem

		movies						
users	*	0.1	*	*	0.9	*	0.9	*
	0.5	0.4	*	*	*	0.6	*	0.6
	*	*	*	0.7	*	*	*	0.2
	0.4	0.2	0.9	*	*	*	*	*
	*	*	0.7	0.7	*	*	0.3	*

$$X^* = \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^{\top}}_{k \times n}$$

$$\boxed{X^*} = \boxed{\tilde{U}} \boxed{\tilde{V}}$$

low-rank

Our problem

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \text{rank}(X) \\ \text{subject to} & X_{ij} = a_{ij}, \quad (i, j) \in \mathcal{O} \end{array}$$

		movies						
users	*	0.1	*	*	0.9	*	0.9	*
	0.5	0.4	*	*	*	0.6	*	0.6
	*	*	*	0.7	*	*	*	0.2
	0.4	0.2	0.9	*	*	*	*	*
	*	*	0.7	0.7	*	*	0.3	*

$$X^* = \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^T}_{k \times n}$$

$$\boxed{X^*} = \boxed{\tilde{U}} \boxed{\tilde{V}}$$

low-rank

Our problem

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \text{rank}(X) \\ \text{subject to} & X_{ij} = a_{ij}, \quad (i, j) \in \mathcal{O} \end{array}$$

observed entries

		movies						
users	*	0.1	*	*	0.9	*	0.9	*
	0.5	0.4	*	*	*	0.6	*	0.6
	*	*	*	0.7	*	*	*	0.2
	0.4	0.2	0.9	*	*	*	*	*
	*	*	0.7	0.7	*	*	0.3	*

$$X^* = \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^{\top}}_{k \times n}$$

$$\boxed{X^*} = \boxed{\tilde{U}} \boxed{\tilde{V}}$$

low-rank

Our problem

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \text{rank}(X) \\ \text{subject to} & X_{ij} = a_{ij}, \quad (i, j) \in \mathcal{O} \end{array}$$

nonconvex

observed entries

		movies						
users	*	0.1	*	*	0.9	*	0.9	*
	0.5	0.4	*	*	*	0.6	*	0.6
	*	*	*	0.7	*	*	*	0.2
	0.4	0.2	0.9	*	*	*	*	*
	*	*	0.7	0.7	*	*	0.3	*

$$X^* = \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^T}_{k \times n}$$

$$\boxed{X^*} = \boxed{\tilde{U}} \boxed{\tilde{V}}$$

low-rank

Our problem

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \text{rank}(X) \\ \text{subject to} & X_{ij} = a_{ij}, \quad (i, j) \in \mathcal{O} \end{array}$$

nonconvex

observed entries

relax

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|(\sigma_1(X), \sigma_2(X), \dots, \sigma_r(X))\|_1 \\ \text{subject to} & X_{ij} = a_{ij}, \quad (i, j) \in \mathcal{O} \end{array}$$

		movies						
users	*	0.1	*	*	0.9	*	0.9	*
	0.5	0.4	*	*	*	0.6	*	0.6
	*	*	*	0.7	*	*	*	0.2
	0.4	0.2	0.9	*	*	*	*	*
	*	*	0.7	0.7	*	*	0.3	*

$$X^* = \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^T}_{k \times n}$$

$$\boxed{X^*} = \boxed{\tilde{U}} \boxed{\tilde{V}}$$

low-rank

Our problem

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \text{rank}(X) \\ \text{subject to} & X_{ij} = a_{ij}, \quad (i, j) \in \mathcal{O} \end{array}$$

nonconvex

observed entries

↓ relax

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|(\sigma_1(X), \sigma_2(X), \dots, \sigma_r(X))\|_1 = \|X\|_* \\ \text{subject to} & X_{ij} = a_{ij}, \quad (i, j) \in \mathcal{O} \end{array}$$

		movies						
users	★	0.1	★	★	0.9	★	0.9	★
	0.5	0.4	★	★	★	0.6	★	0.6
	★	★	★	0.7	★	★	★	0.2
	0.4	0.2	0.9	★	★	★	★	★
	★	★	0.7	0.7	★	★	0.3	★

$$X^* = \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^T}_{k \times n}$$

$$\boxed{X^*} = \boxed{\tilde{U}} \boxed{\tilde{V}}$$

low-rank

Our problem

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \text{rank}(X) \\ \text{subject to} & X_{ij} = a_{ij}, \quad (i, j) \in \mathcal{O} \end{array}$$

nonconvex

observed entries

relax

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|(\sigma_1(X), \sigma_2(X), \dots, \sigma_r(X))\|_1 \\ \text{subject to} & X_{ij} = a_{ij}, \quad (i, j) \in \mathcal{O} \end{array}$$

nuclear norm

$$= \|X\|_{\star}$$

		movies						
users	★	0.1	★	★	0.9	★	0.9	★
	0.5	0.4	★	★	★	0.6	★	0.6
	★	★	★	0.7	★	★	★	0.2
	0.4	0.2	0.9	★	★	★	★	★
	★	★	0.7	0.7	★	★	0.3	★

$$X^* = \underbrace{U \Sigma^{\frac{1}{2}}}_{m \times k} \underbrace{\Sigma^{\frac{1}{2}} V^T}_{k \times n}$$

$$\boxed{X^*} = \boxed{\tilde{U}} \boxed{\tilde{V}}$$

low-rank

Our problem

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \text{rank}(X) \\ \text{subject to} & X_{ij} = a_{ij}, \quad (i, j) \in \mathcal{O} \end{array}$$

nonconvex

observed entries

↓ relax

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|(\sigma_1(X), \sigma_2(X), \dots, \sigma_r(X))\|_1 \\ \text{subject to} & X_{ij} = a_{ij}, \quad (i, j) \in \mathcal{O} \end{array}$$

nuclear norm

$$= \|X\|_{\star}$$

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|X\|_{\star} \\ \text{subject to} & \text{tr}(X M_l) = a_l, \quad l = 1, \dots, p \end{array}$$

		movies							
users	★	0.1	★	★	0.9	★	0.9	★	
	0.5	0.4	★	★	★	0.6	★	0.6	
	★	★	★	0.7	★	★	★	0.2	
	0.4	0.2	0.9	★	★	★	★	★	
	★	★	0.7	0.7	★	★	0.3	★	

Example Result

Example Result

Theorem [Chandrasekaran et al. 12']

Example Result

Theorem [Chandrasekaran et al. 12']

$$X^* \in \mathbb{R}^{m \times n}$$

Example Result

Theorem [Chandrasekaran et al. 12']

$X^* \in \mathbb{R}^{m \times n}$ *unknown*, but rank k

Example Result

Theorem [Chandrasekaran et al. 12']

$X^* \in \mathbb{R}^{m \times n}$ *unknown*, but rank k

$a_l = \text{tr}(X M_l)$, $l = 1, \dots, p$ *measurements*

Example Result

Theorem [Chandrasekaran et al. 12']

$X^* \in \mathbb{R}^{m \times n}$ *unknown*, but rank k

$a_l = \text{tr}(X M_l)$, $l = 1, \dots, p$ *measurements*
| iid entries $\mathcal{N}(0, 1)$

Example Result

Theorem [Chandrasekaran et al. 12']

$X^* \in \mathbb{R}^{m \times n}$ *unknown*, but rank k

$$p \geq 3k(m + n - k) + 1$$

iid entries $\mathcal{N}(0, 1)$
 $a_l = \text{tr}(X M_l), \quad l = 1, \dots, p$ *measurements*

Example Result

Theorem [Chandrasekaran et al. 12']

$X^* \in \mathbb{R}^{m \times n}$ *unknown*, but rank k

$$p \geq 3k(m + n - k) + 1$$

\implies

$$a_l = \text{tr}(X M_l), \quad l = 1, \dots, p \quad \text{measurements}$$

| iid entries $\mathcal{N}(0, 1)$

$$X^* = \underset{X}{\text{argmin}} \quad \|X\|_* \quad \text{w.h.p.}$$

s.t. $\text{tr}(X M_l) = a_l, \quad l = 1, \dots, p$

Experiments

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|X\|_{\star} \\ \text{subject to} & \text{tr}(X M_l) = a_l, \quad l = 1, \dots, p \end{array}$$

Experiments

$X^* : 30 \times 30$

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|X\|_* \\ \text{subject to} & \text{tr}(X M_l) = a_l, \quad l = 1, \dots, p \end{array}$$

Experiments

$$X^* : 30 \times 30 \quad \text{rank}(X^*) = 3$$

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|X\|_* \\ \text{subject to} & \text{tr}(X M_l) = a_l, \quad l = 1, \dots, p \end{array}$$

Experiments

$$X^* : 30 \times 30$$

$$\text{rank}(X^*) = 3$$

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|X\|_* \\ \text{subject to} & \text{tr}(X M_l) = a_l, \quad l = 1, \dots, p \end{array}$$

| iid entries $\mathcal{N}(0, 1)$

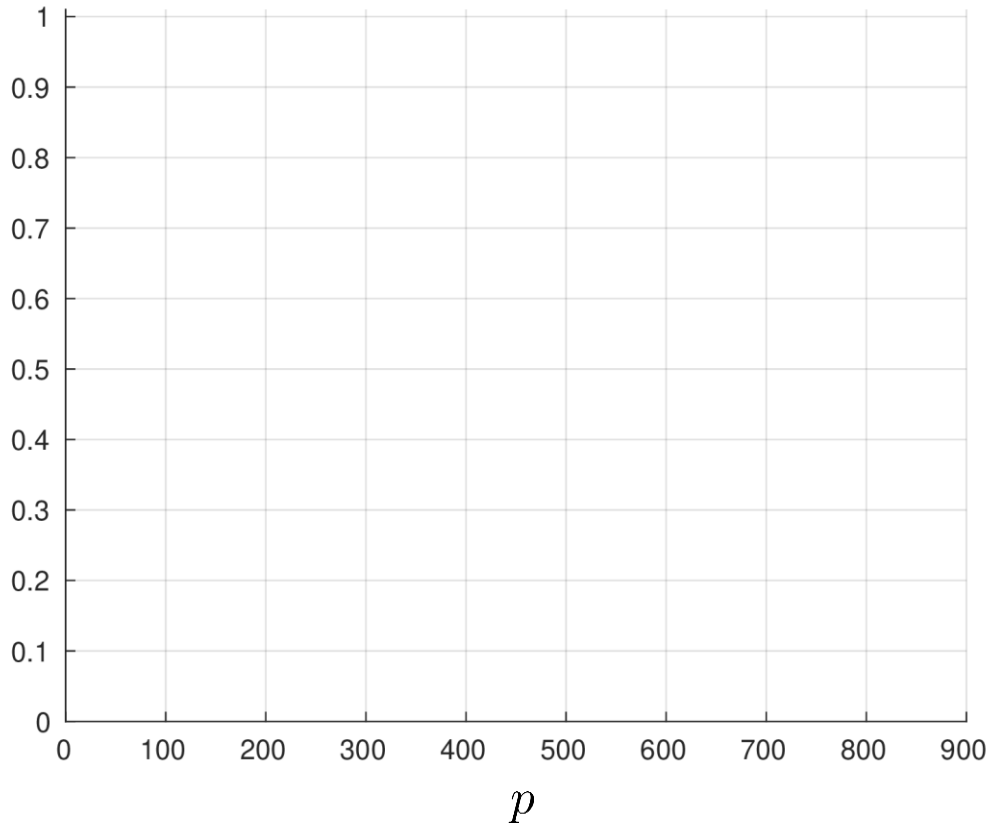
Experiments

$$X^* : 30 \times 30 \quad \text{rank}(X^*) = 3$$

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|X\|_* \\ \text{subject to} & \text{tr}(X M_l) = a_l, \quad l = 1, \dots, p \end{array}$$

| iid entries $\mathcal{N}(0, 1)$

Success rate (20 trials)



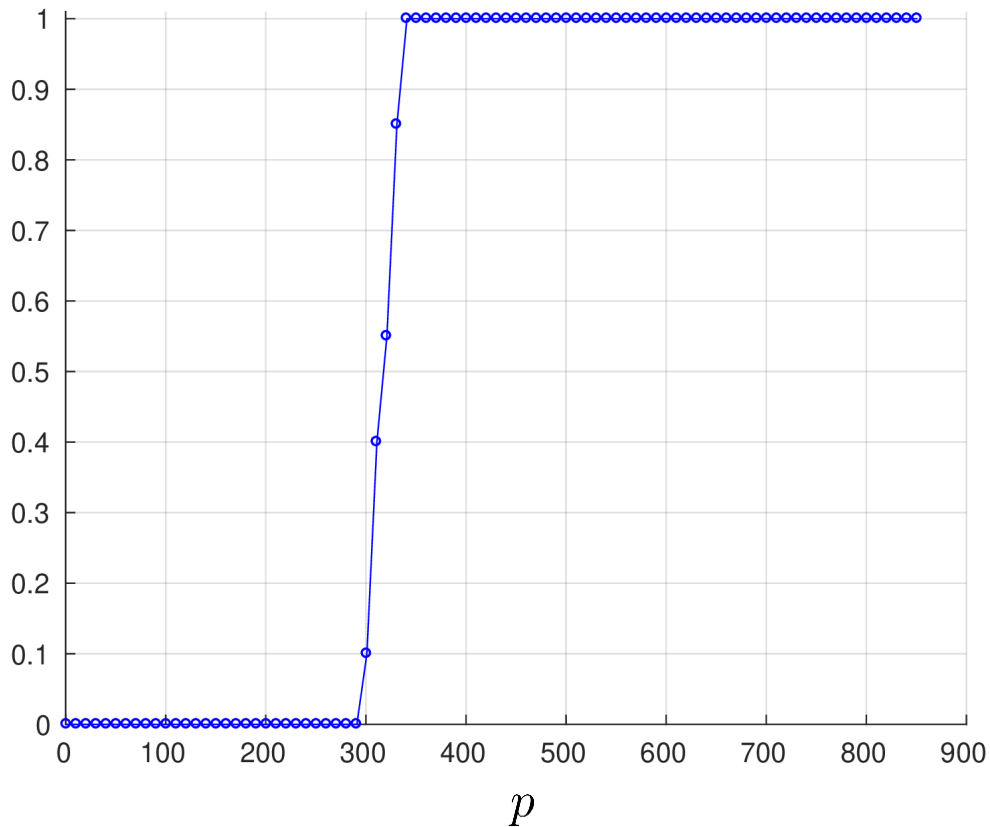
Experiments

 $X^* : 30 \times 30$ $\text{rank}(X^*) = 3$

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|X\|_* \\ \text{subject to} & \text{tr}(X M_l) = a_l, \quad l = 1, \dots, p \end{array}$$

| iid entries $\mathcal{N}(0, 1)$

Success rate (20 trials)



Experiments

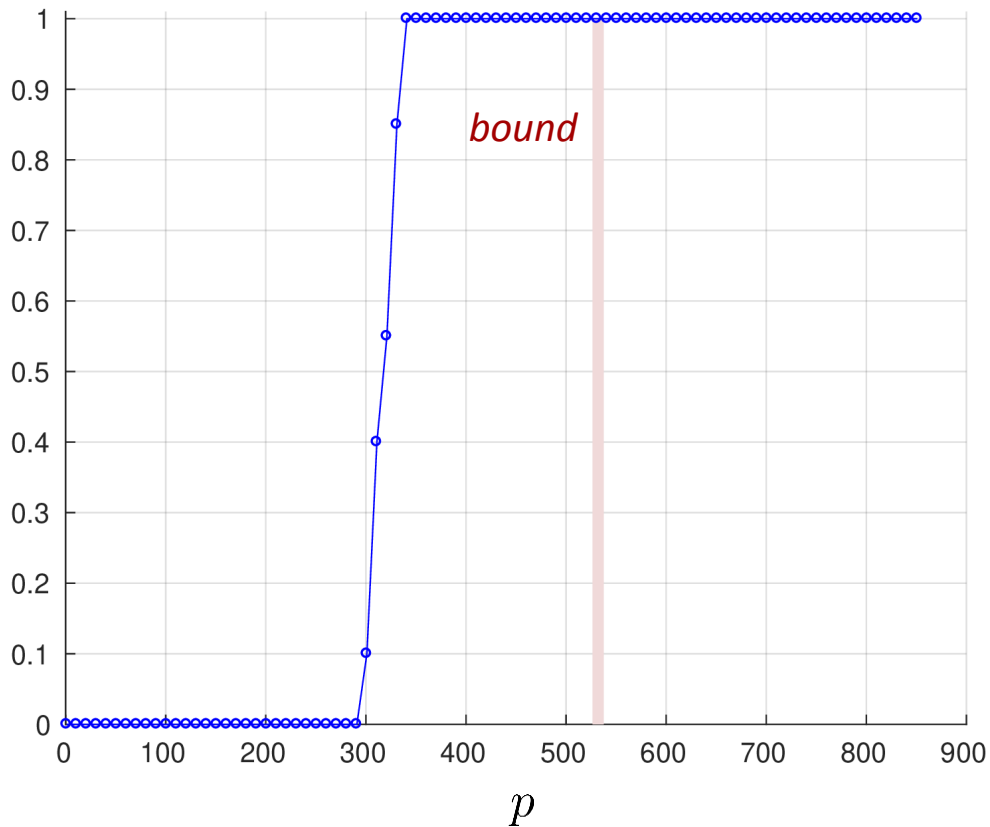
$X^* : 30 \times 30$

$\text{rank}(X^*) = 3$

$$\begin{array}{ll} \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} & \|X\|_* \\ \text{subject to} & \text{tr}(X M_l) = a_l, \quad l = 1, \dots, p \end{array}$$

iid entries $\mathcal{N}(0, 1)$

Success rate (20 trials)



Experiments

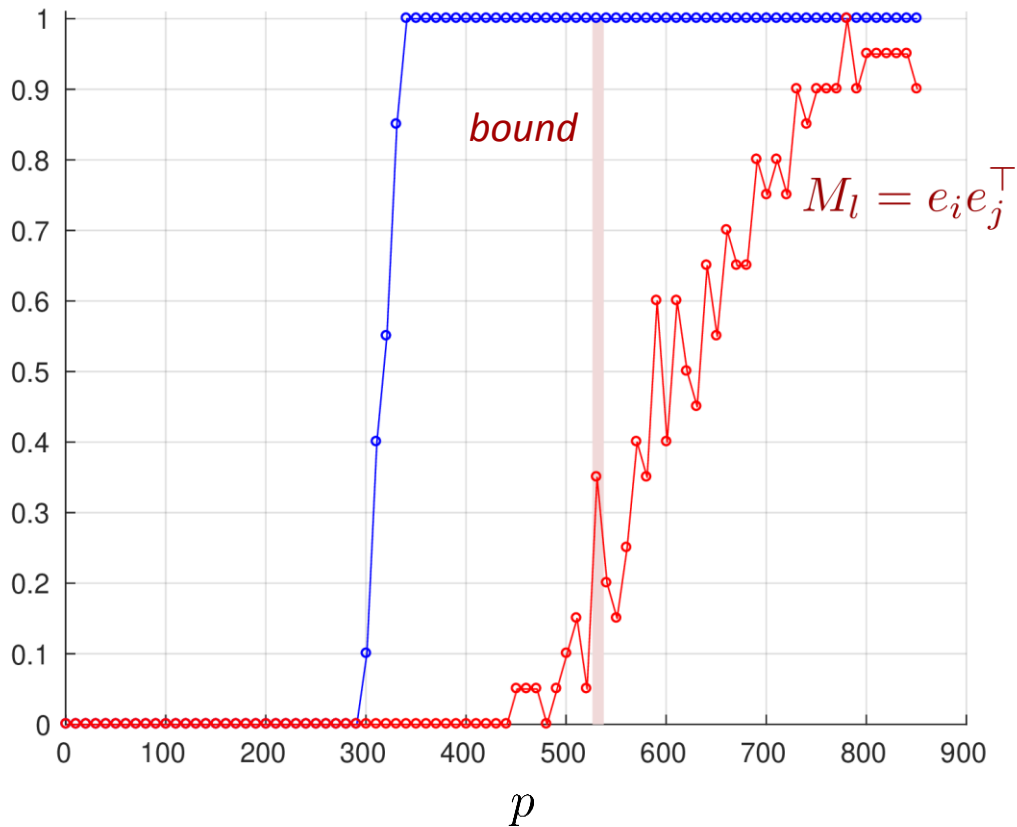
$X^* : 30 \times 30$

$\text{rank}(X^*) = 3$

$$\begin{aligned} & \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} && \|X\|_* \\ & \text{subject to} && \text{tr}(X M_l) = a_l, \quad l = 1, \dots, p \end{aligned}$$

iid entries $\mathcal{N}(0, 1)$

Success rate (20 trials)



Experiments

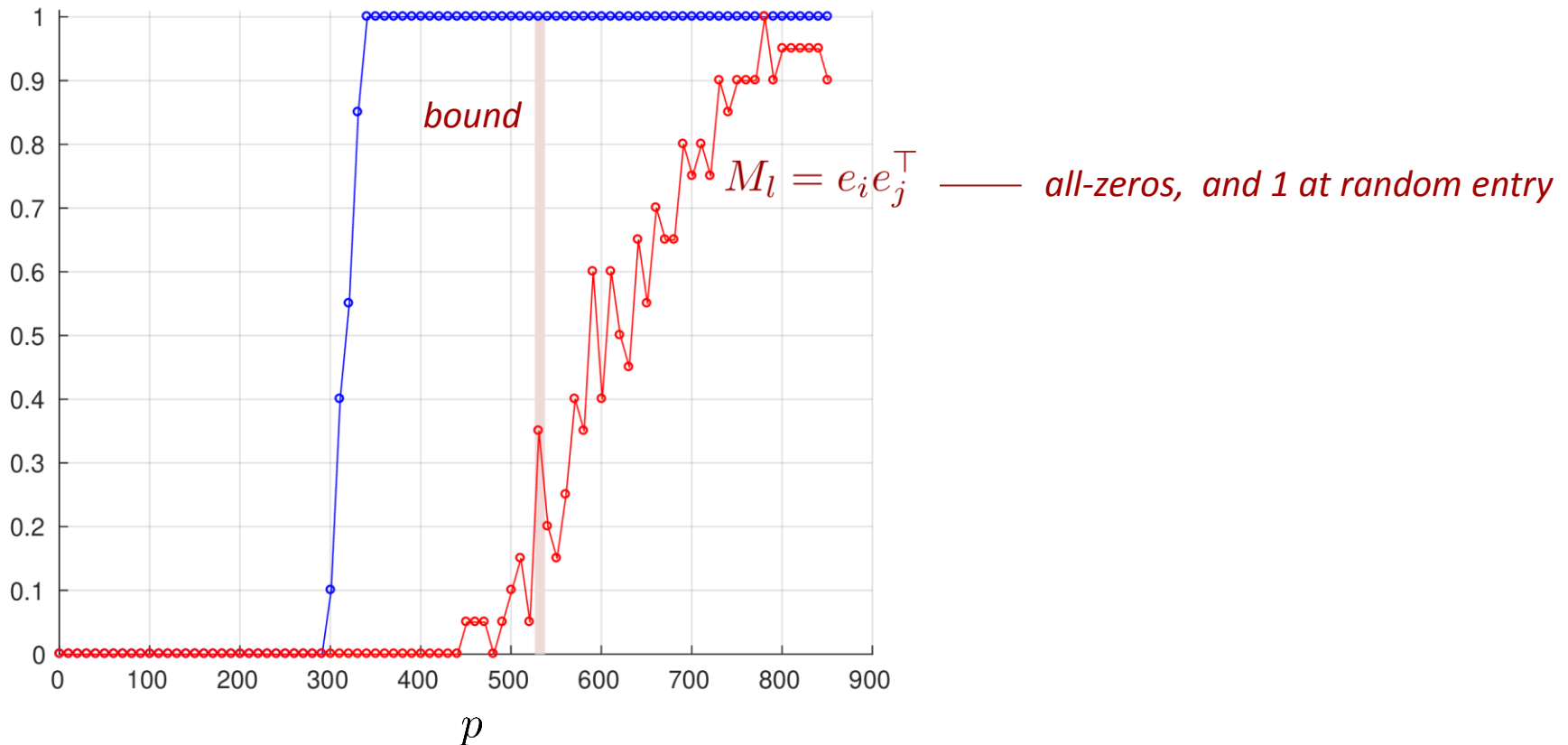
$X^* : 30 \times 30$

$\text{rank}(X^*) = 3$

$$\begin{aligned} &\text{minimize}_{X \in \mathbb{R}^{m \times n}} && \|X\|_* \\ &\text{subject to} && \text{tr}(X M_l) = a_l, \quad l = 1, \dots, p \end{aligned}$$

iid entries $\mathcal{N}(0, 1)$

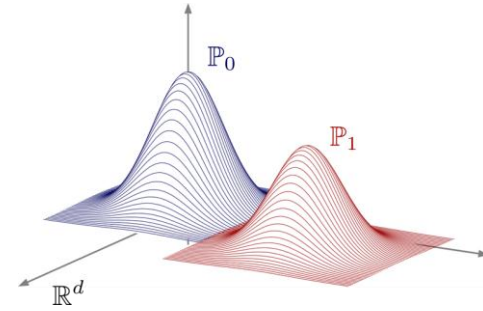
Success rate (20 trials)



Conclusions

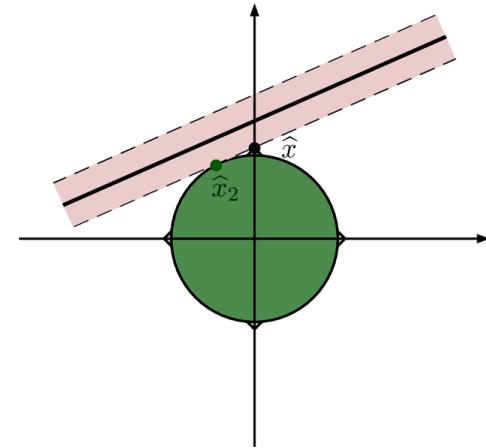
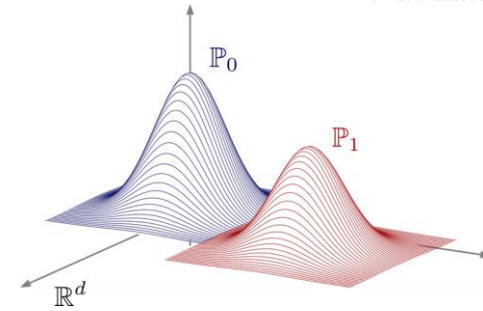
Conclusions

- Structure is key in *high-dimensional* problems



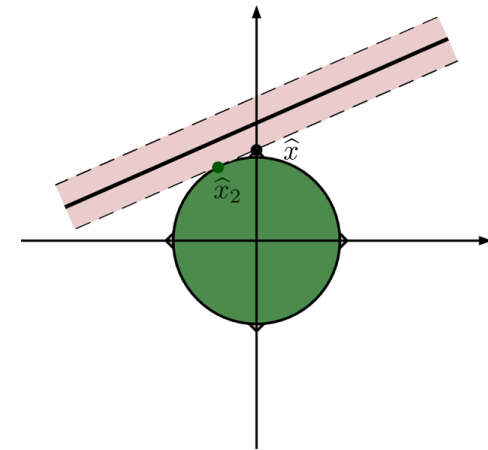
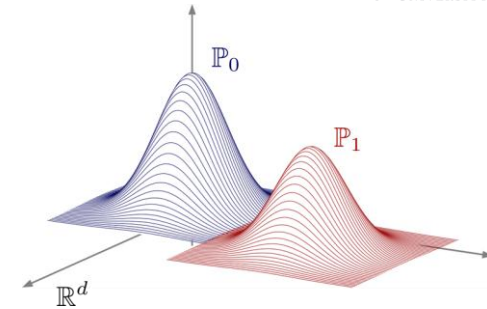
Conclusions

- Structure is key in *high-dimensional* problems
- *Sparsity* encodes several types of structure

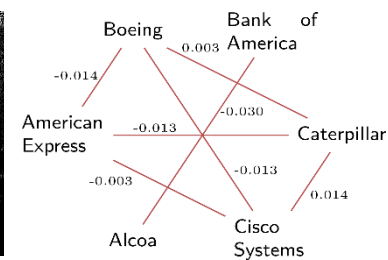
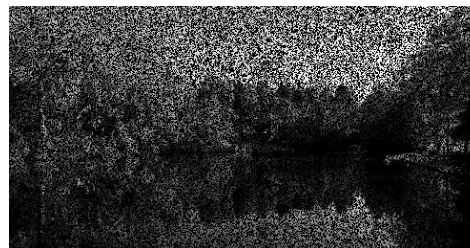


Conclusions

- Structure is key in *high-dimensional* problems
- *Sparsity* encodes several types of structure
- Several applications (and theory)

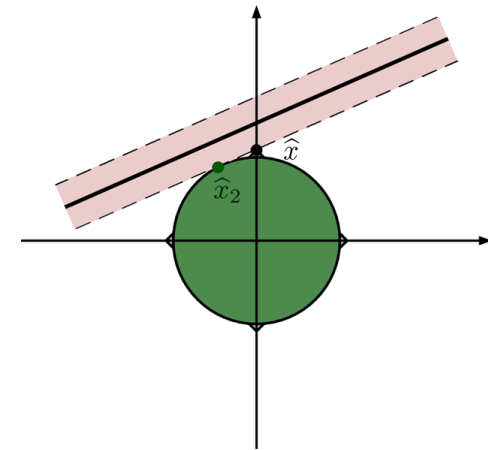
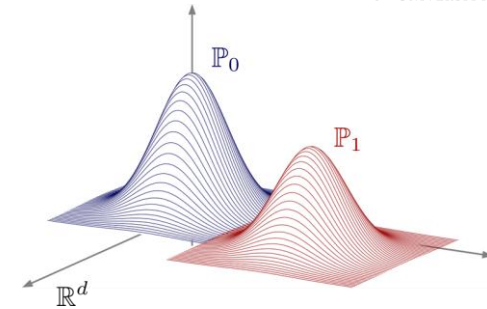


City	Police force/1000 (1990)	% with 4 years high school (1990)	% not in high school (15-19)	% in college (18-24)	% with 4+ years college (25+)	Crime (millions)
1	40	74	11	31	20	478
2	32	72	11	43	18	404
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
...
50	66	67	26	18	16	940

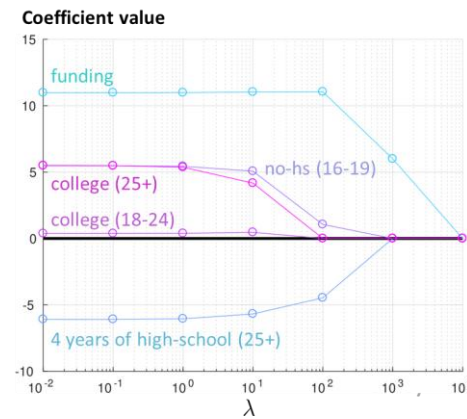
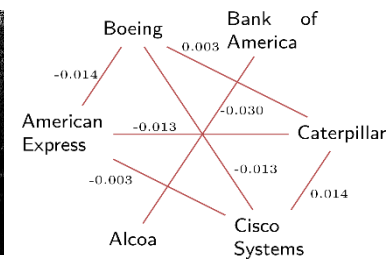


Conclusions

- Structure is key in *high-dimensional* problems
- *Sparsity* encodes several types of structure
- Several applications (and theory)
- LASSO, basis pursuit, ... improve *interpretability* and (often) *performance*

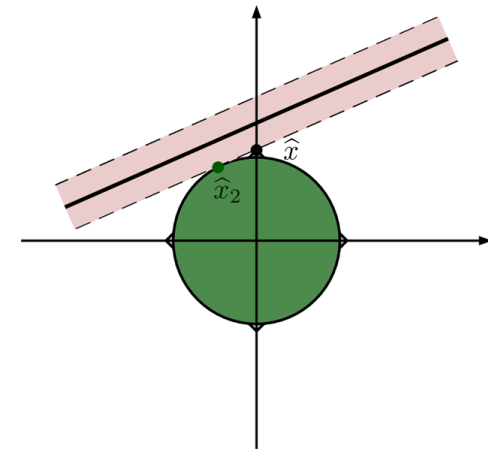
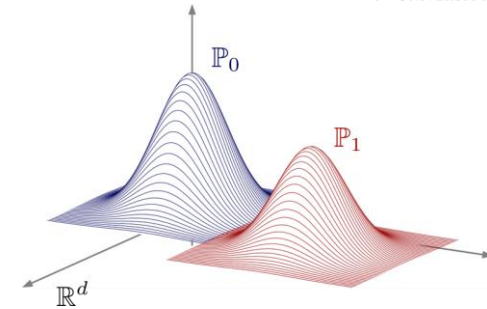


City	Police force per resident (1/100)	% with 4 years high-school (15-25)	% not in high school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Census (million)
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
...
50	66	67	26	18	16	940

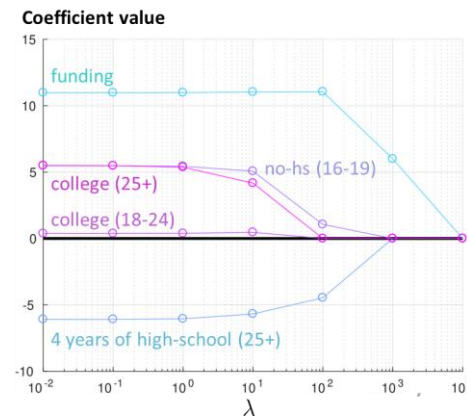
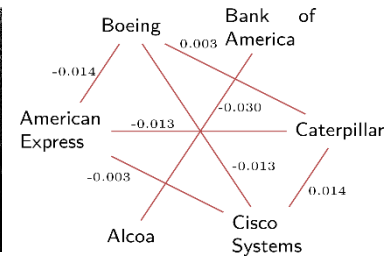
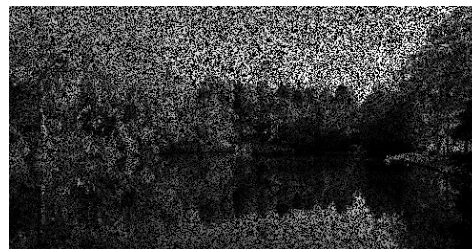


Conclusions

- Structure is key in *high-dimensional* problems
- *Sparsity* encodes several types of structure
- Several applications (and theory)
- LASSO, basis pursuit, ... improve *interpretability* and (often) *performance*
- Didn't cover: optimization theory and *algorithms*



City	Police force per resident (1/100)	% with 4 years high-school (15-25)	% not in high school (16-19)	% in college (18-24)	% with 4+ years college (25+)	Census (million)
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
6	25	68	8	32	15	603
7	34	68	12	24	14	484
8	33	62	13	28	11	546
9	36	69	7	25	12	424
...
50	66	67	26	18	16	940



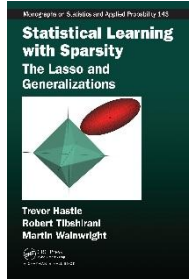
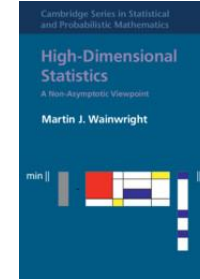
References

References

M. J. Wainwright

High-Dimensional Statistics: A Non-Asymptotic Viewpoint

Cambridge University Press, 2019



T. Hastie, R. Tibshirani, M. Wainwright

Statistical Learning with Sparsity: The Lasso and Generalizations

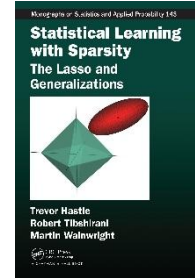
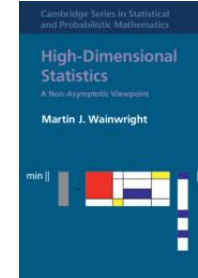
CRC Press, 2016 <https://web.stanford.edu/~hastie/StatLearnSparsity/>

References

M. J. Wainwright

High-Dimensional Statistics: A Non-Asymptotic Viewpoint

Cambridge University Press, 2019



T. Hastie, R. Tibshirani, M. Wainwright

Statistical Learning with Sparsity: The Lasso and Generalizations

CRC Press, 2016 <https://web.stanford.edu/~hastie/StatLearnSparsity/>

V. Chandrasekaran, B. Recht, P. A. Parrilo, A. S. Willsky

The Convex Geometry of Linear Inverse Problems

Foundations of Computational Mathematics, Vol. 12, pp. 805-849, 2012

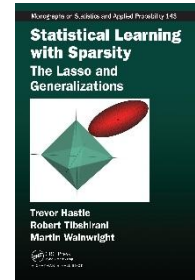
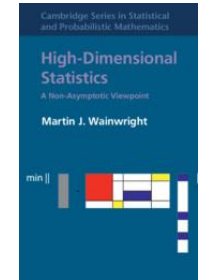


References

M. J. Wainwright

High-Dimensional Statistics: A Non-Asymptotic Viewpoint

Cambridge University Press, 2019



T. Hastie, R. Tibshirani, M. Wainwright

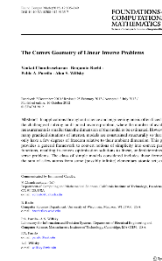
Statistical Learning with Sparsity: The Lasso and Generalizations

CRC Press, 2016 <https://web.stanford.edu/~hastie/StatLearnSparsity/>

V. Chandrasekaran, B. Recht, P. A. Parrilo, A. S. Willsky

The Convex Geometry of Linear Inverse Problems

Foundations of Computational Mathematics, Vol. 12, pp. 805-849, 2012



J. F. C. Mota, N. Deligiannis, M. R. D. Rodrigues

Compressed sensing with side information: Geometrical interpretation and performance bounds

IEEE Global Conference on Signal and Information Processing (GlobalSIP), 2014

P. Ravikumar, M. J. Wainwright, G. Raskutti, B. Yu

High-dimensional covariance estimation by minimizing L1-penalized log-determinant divergence

Electronic Journal of Statistics, Vol. 5, pp. 935-980, 2011

References

V. Chandrasekaran

Convex Optimization Methods for Graphs and Statistical Modeling

PhD thesis, MIT, 2011

M. S. Brown, M. Pelosi, H. Dirska

Dynamic-radius species-conserving genetic algorithm for the financial forecasting of Dow Jones index stocks

Machine Learning and Data Mining in Pattern Recognition, Vol. 7988, pp. 27-41, 2013

Code & presentation

<https://github.com/joaofcmota/udrc-summerschool>

<http://jmota.eps.hw.ac.uk/documents/Mota21-HighDimensionalStatsAndSparsity-UDRC.pdf>