

# Model Selection in High Dimension

---

Jalal Fadili

Normandie Université-ENSICAEN, GREYC

*Mathematical coffees 2017*



Normandie Université



ÉCOLE PUBLIQUE D'INGÉNIEURS  
CENTRE DE RECHERCHE

# Model selection: what for ?

- A key conceptual tool for **dimension reduction** and exploiting hidden structures in **high-dimensional** data.



# Model selection: what for ?

- A key conceptual tool for **dimension reduction** and exploiting hidden structures in **high-dimensional** data.
- The gist:
  - **Compare** different statistical models each for a possible hidden structure.
  - **Select** the one that is more suited for your task.



# Model selection: what for ?

- A key conceptual tool for **dimension reduction** and exploiting hidden structures in **high-dimensional** data.
- The gist:
  - **Compare** different statistical models each for a possible hidden structure.
  - **Select** the one that is more suited for your task.
- Challenges and questions:
  - Data-driven model selection, i.e. without an oracle.
  - Guarantees ? In terms of what ?
  - Optimality.
  - Computational issues.
  - Tractable procedures (recall MC on sparsity and CS).

# Gaussian regression

$$y_i = f^*(x_i) + \varepsilon_i$$

- $x_1, \dots, x_n$  : design vectors in  $\mathbb{R}^p$ .
- $f^* : \mathbb{R}^p \rightarrow \mathbb{R}$  : unknown regression function.
- $(\varepsilon_1, \dots, \varepsilon_n)$  are independent and identically distributed  $\mathcal{N}(0, \sigma^2)$ .

# Gaussian regression

$$y_i = f^*(x_i) + \varepsilon_i$$

- $x_1, \dots, x_n$  : design vectors in  $\mathbb{R}^p$ .
- $f^* : \mathbb{R}^p \rightarrow \mathbb{R}$  : unknown regression function.
- $(\varepsilon_1, \dots, \varepsilon_n)$  are independent and identically distributed  $\mathcal{N}(0, \sigma^2)$ .

## Goal

Estimate  $f^*$  from data  $\{(x_i, y_i)\}_{1 \leq i \leq n}$

# Gaussian regression

$$y_i = f^*(x_i) + \varepsilon_i$$

- $x_1, \dots, x_n$  : design vectors in  $\mathbb{R}^p$ .
- $f^* : \mathbb{R}^p \rightarrow \mathbb{R}$  : unknown regression function.
- $(\varepsilon_1, \dots, \varepsilon_n)$  are independent and identically distributed  $\mathcal{N}(0, \sigma^2)$ .
- Curse of dimensionality :
  - $f^*$  Lipschitz.
  - Need  $n \asymp (1/\sigma)^{p+2}$  for root mean-square error  $\sigma$ .
- Blessings of dimensionality :
  - Concentration of measure.
  - High-dimensional geometry of convex bodies.
  - Asymptotic regime.
  - Approach to continuum.

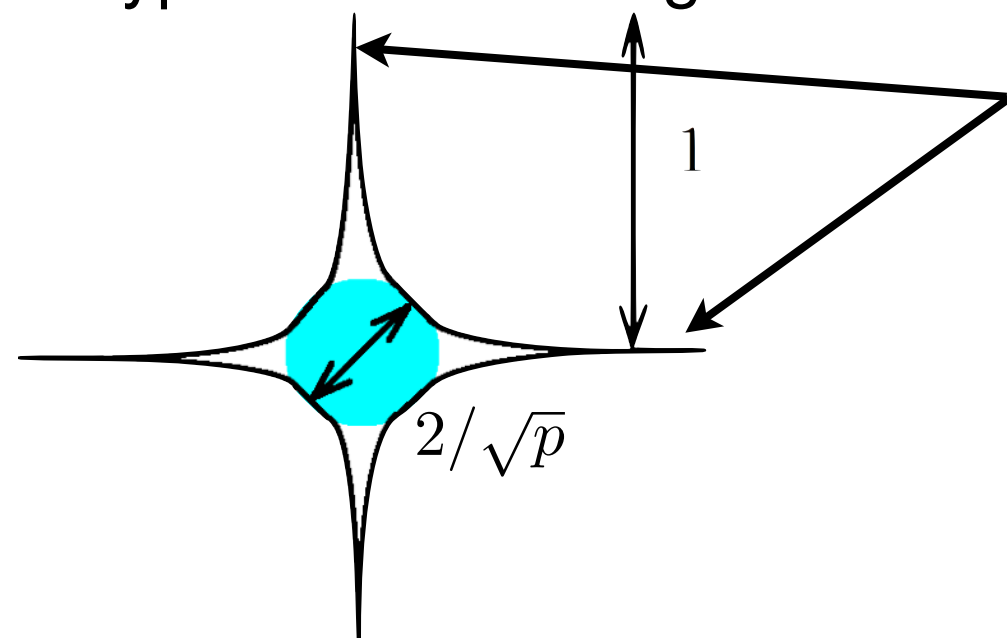
# Convex bodies in high-dimension

- A convex body  $C$  usually consists of a bulk and outliers.
- The bulk makes up most of the volume of  $C$ , but it is usually small in diameter.
- The outliers contribute little to the volume, but they are large in diameter.
- If  $C$  is properly scaled, the bulk usually looks like a Euclidean ball.
- The outliers look like thin, long tentacles.
- The volume in high dimensions scales differently than in low dimensions :
  - dilate by 2 increases volume by  $2^p$ .
  - the tentacles contain exponentially less volume than the bulk.
- This is better captured in a "hyperbolic" drawing.



# Convex bodies in high-dimension

- A convex body  $C$  usually consists of a bulk and outliers.
- The bulk makes up most of the volume of  $C$ , but it is usually small in diameter.
- The outliers contribute little to the volume, but they are large in diameter.
- If  $C$  is properly scaled, the bulk usually looks like a Euclidean ball.
- The outliers look like thin, long tentacles.
- The volume in high dimensions scales differently than in low dimensions :
  - dilate by 2 increases volume by  $2^p$ .
  - the tentacles contain exponentially less volume than the bulk.
- This is better captured in a "hyperbolic" drawing.



**Tentacles protrude quite far along the coordinate directions**

**Milman's "hyperbolic" drawing of the unit  $\ell_1$ -ball**

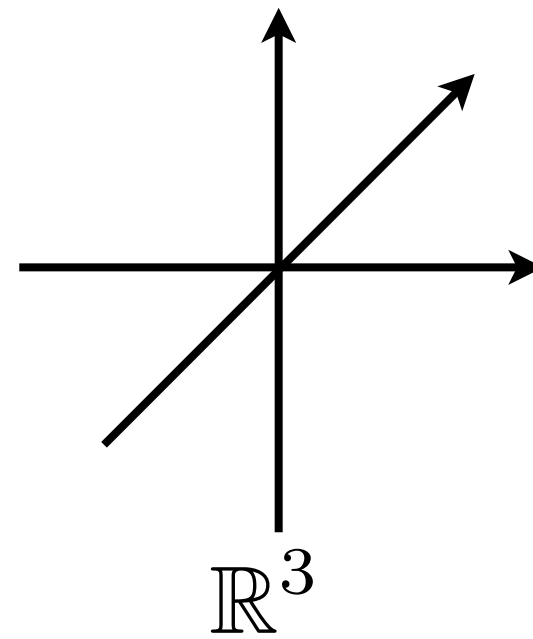
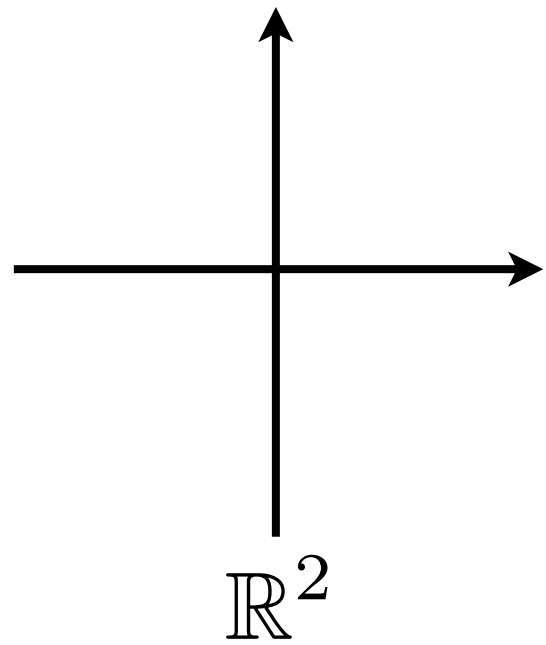
# Blessings in action

$$y \in \mathbb{R}^n = X \in \mathbb{R}^{n \times p} \beta^* \in \mathbb{R}^p + \varepsilon \in \mathbb{R}^n$$

Estimate  $\beta^*$  from data  $(X, y)$

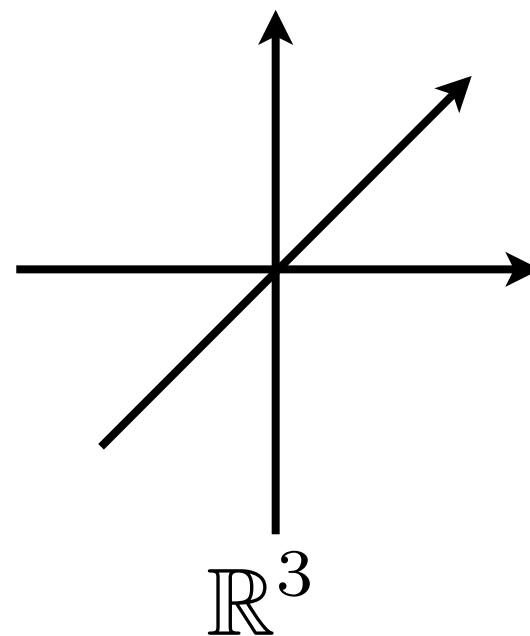
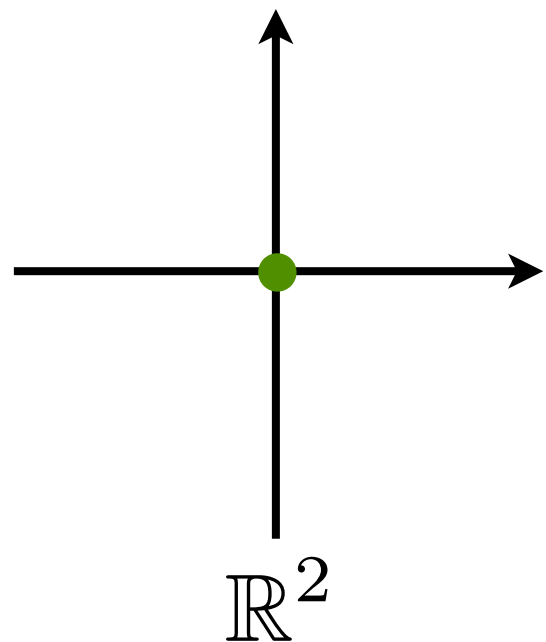
- $X$  is a fat design matrix :
- **High-dimensional** setting : tendency to large  $p$ .
- $\beta^*$  has some **intrinsically small-dimensional** structure :
  - Presumably a few key attributes.
  - Unwilling to specify in advance.
- $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$ .

# Blessings in action: Sparsity



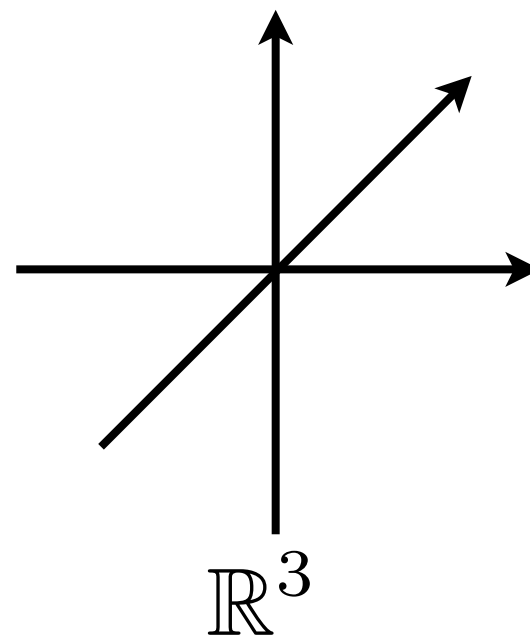
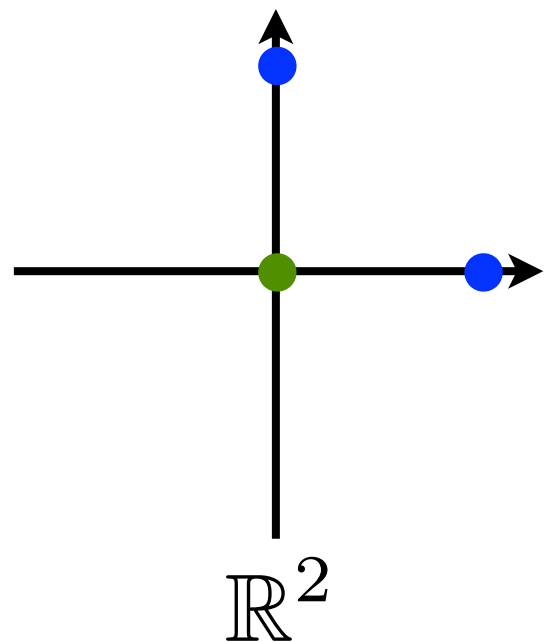
# Blessings in action: Sparsity

- 0-sparse



# Blessings in action: Sparsity

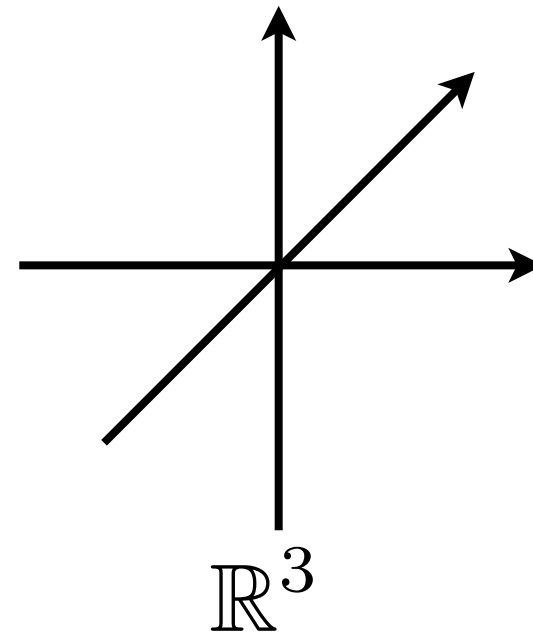
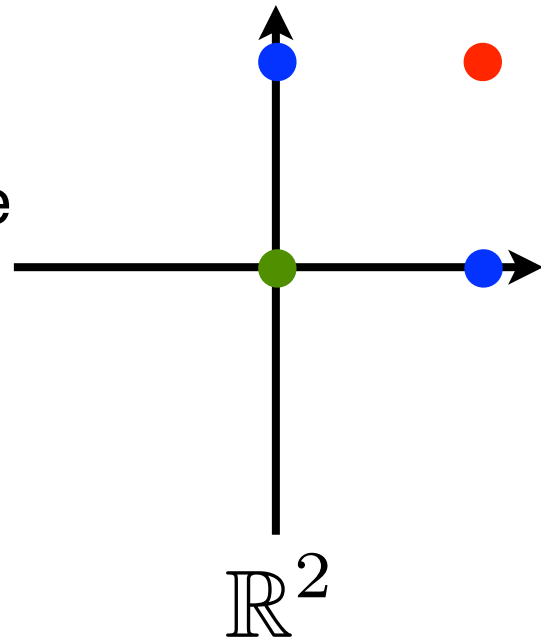
- 0-sparse
- 1-sparse





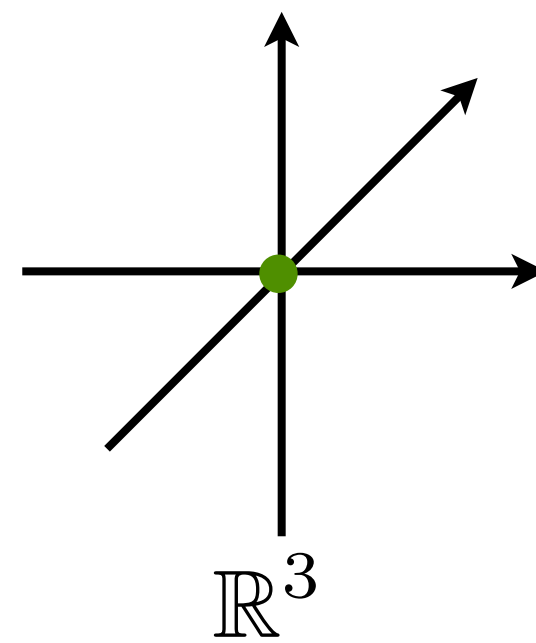
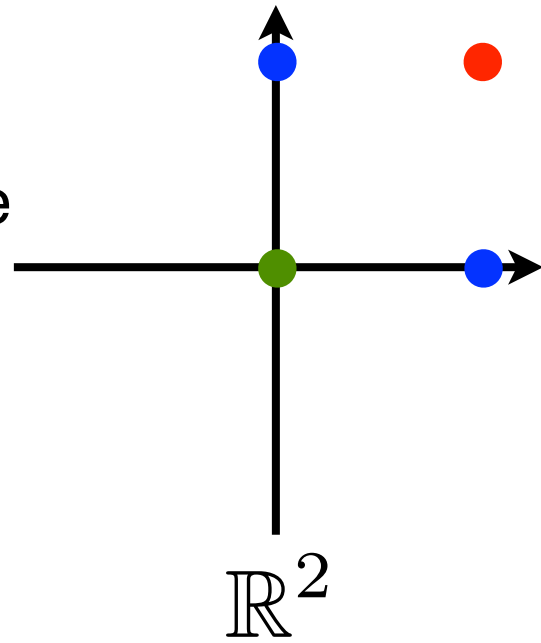
# Blessings in action: Sparsity

- 0-sparse
- 1-sparse
- 2-sparse = dense



# Blessings in action: Sparsity

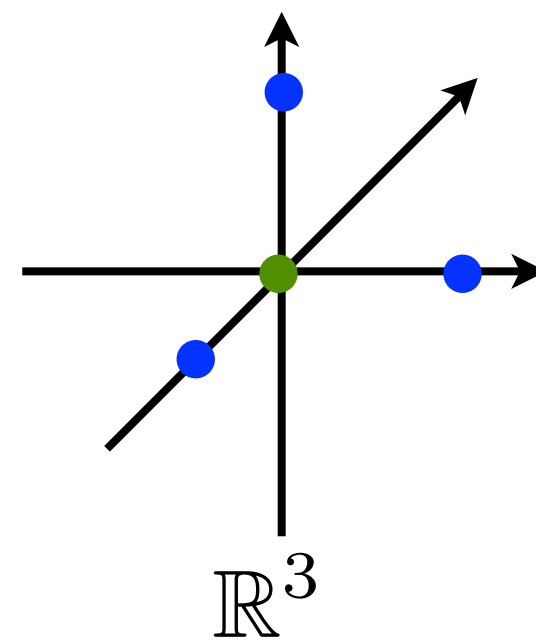
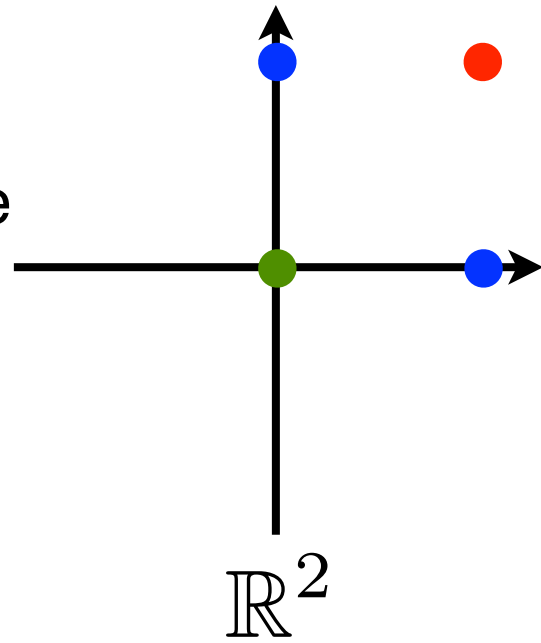
- 0-sparse
- 1-sparse
- 2-sparse = dense



- 0-sparse

# Blessings in action: Sparsity

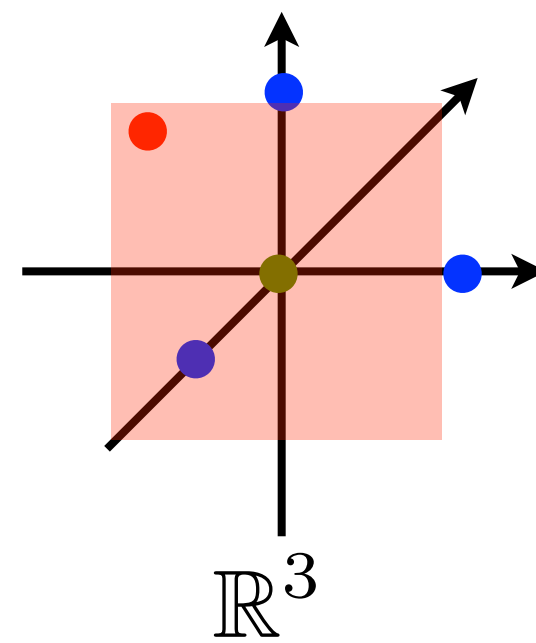
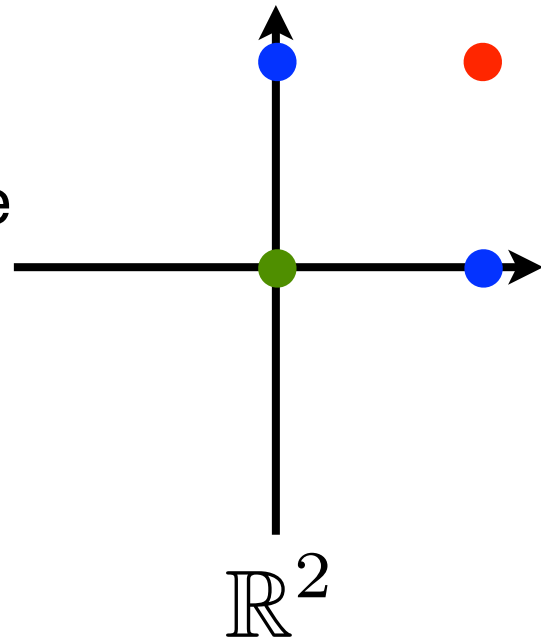
- 0-sparse
- 1-sparse
- 2-sparse = dense



- 0-sparse
- 1-sparse

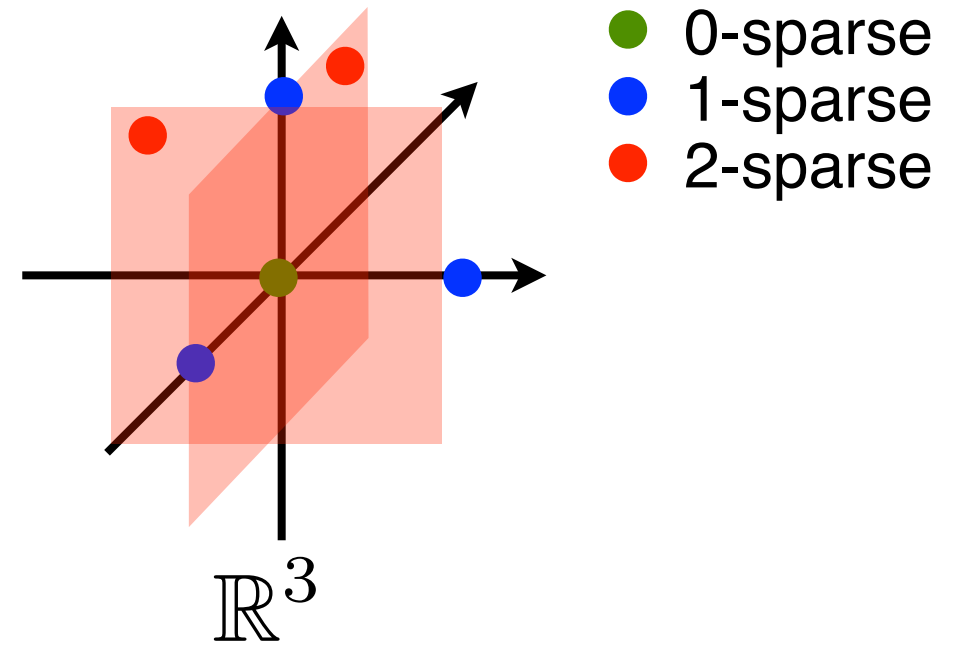
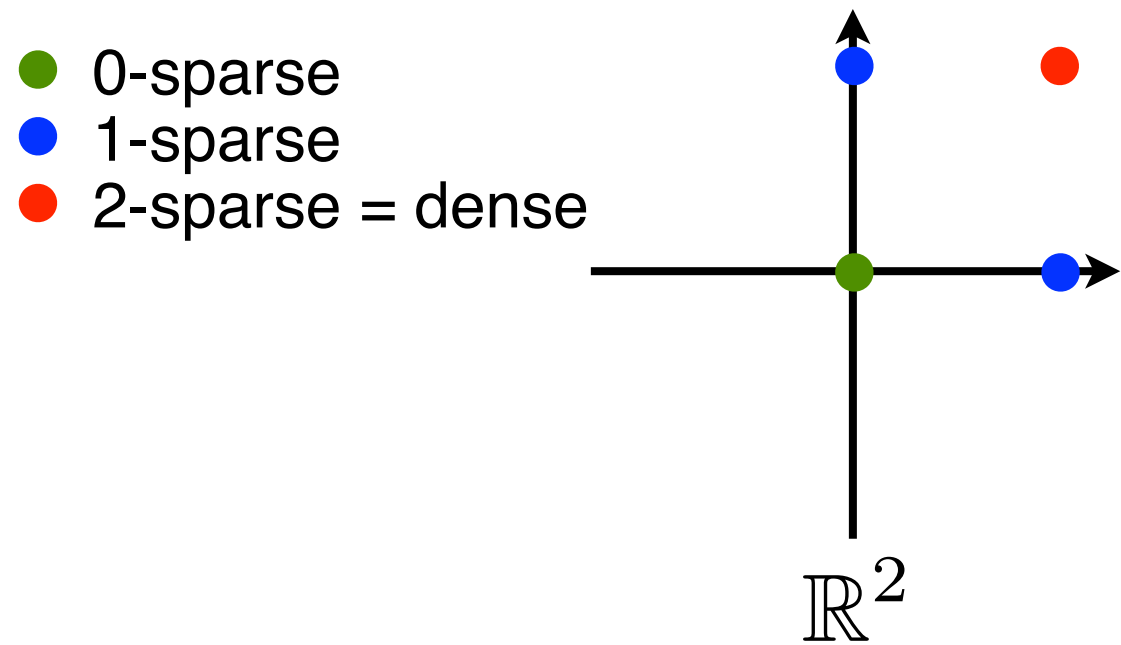
# Blessings in action: Sparsity

- 0-sparse
- 1-sparse
- 2-sparse = dense



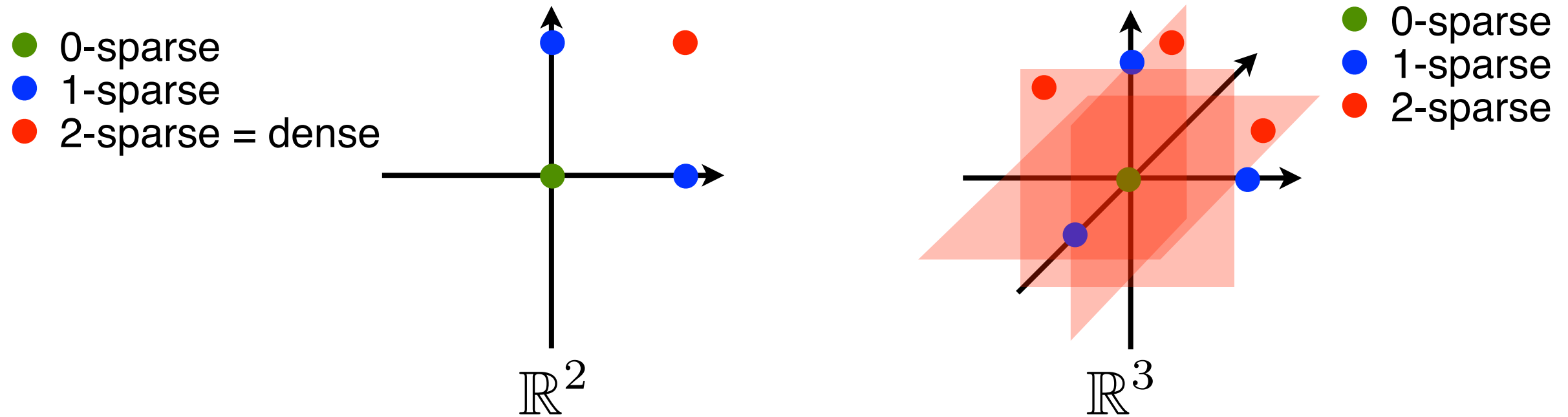
- 0-sparse
- 1-sparse
- 2-sparse

# Blessings in action: Sparsity

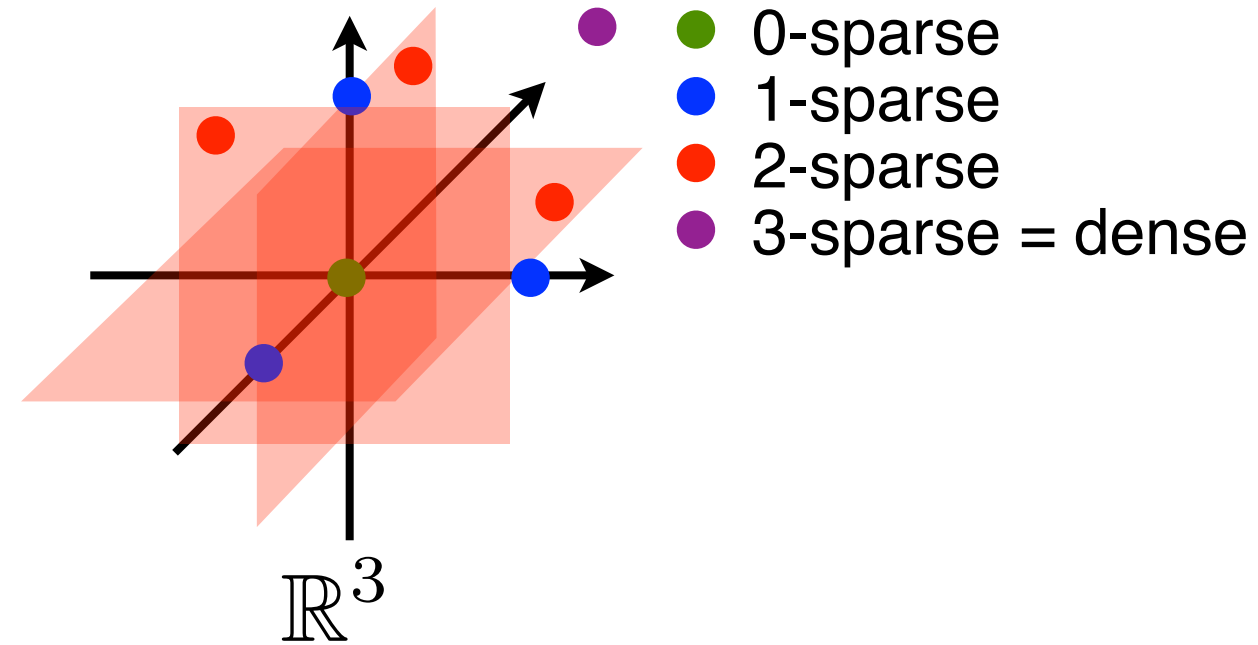
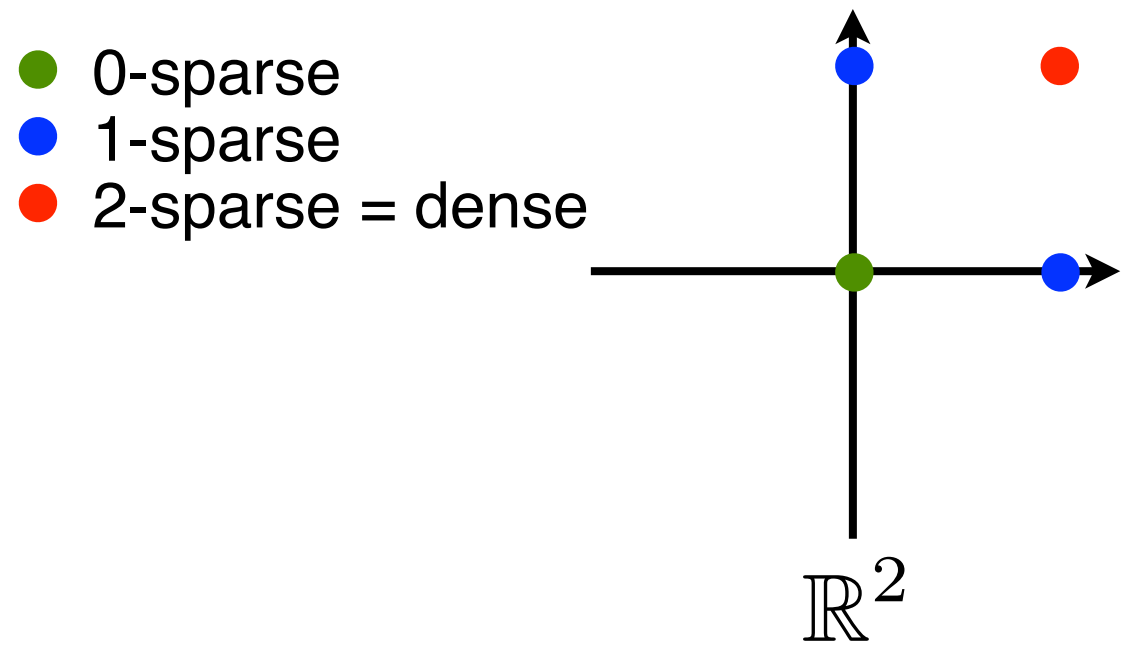




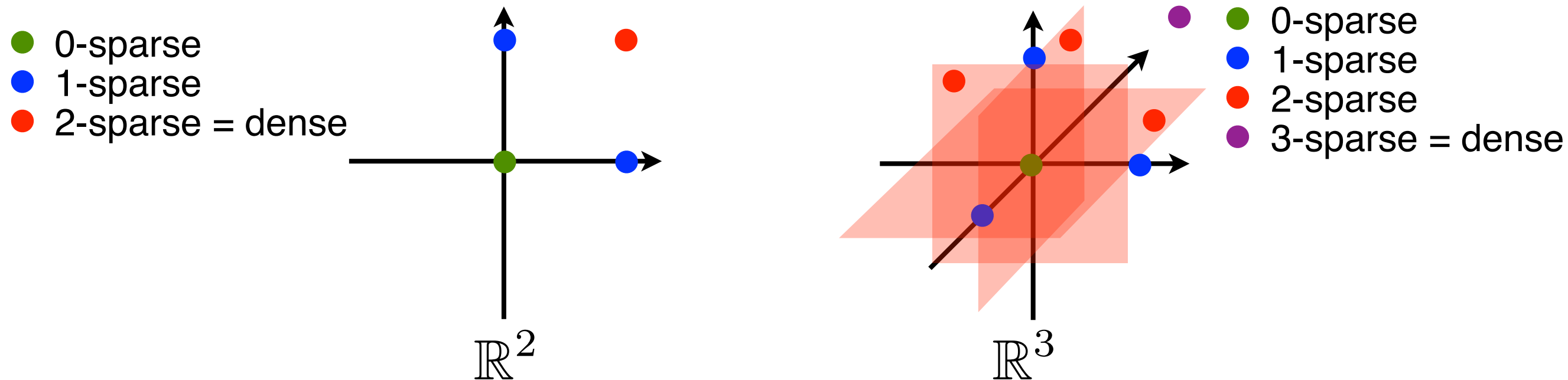
# Blessings in action: Sparsity



# Blessings in action: Sparsity



# Blessings in action: Sparsity

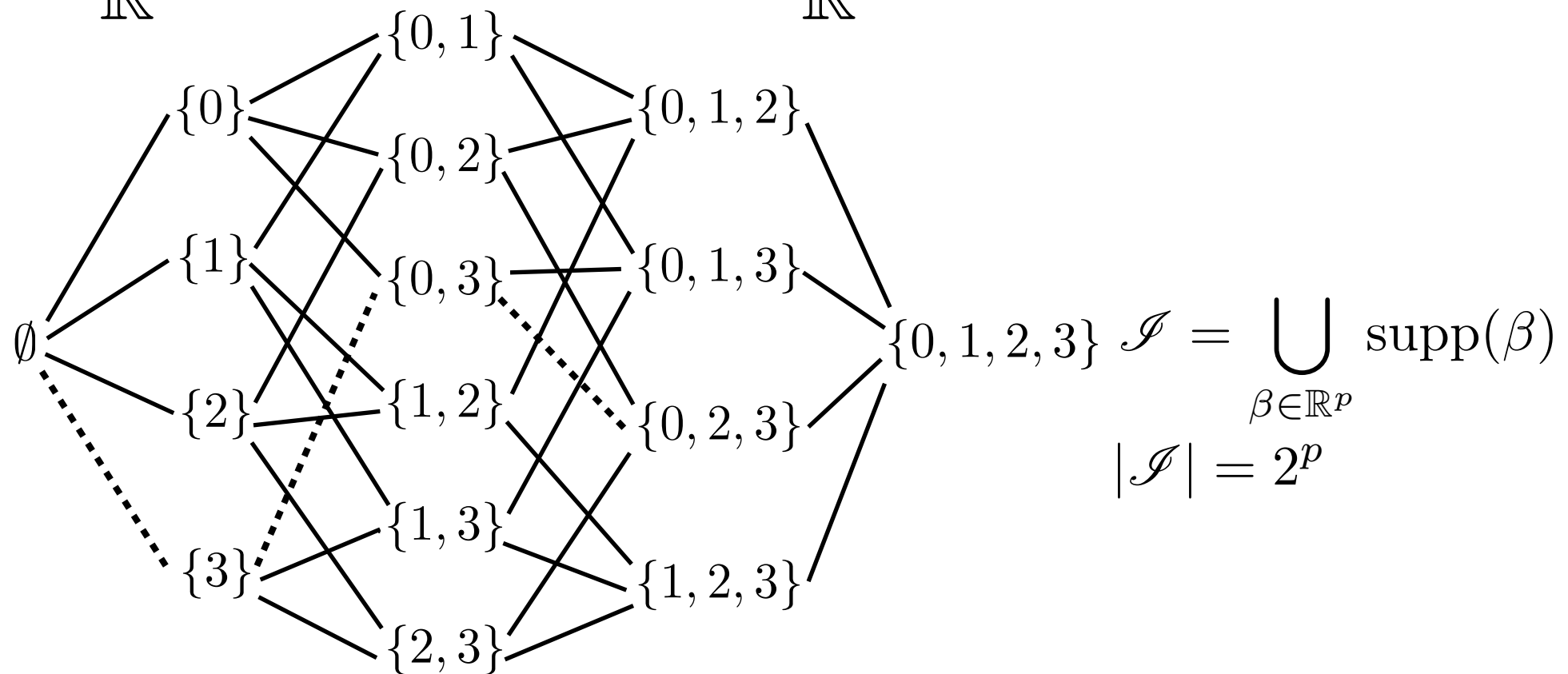
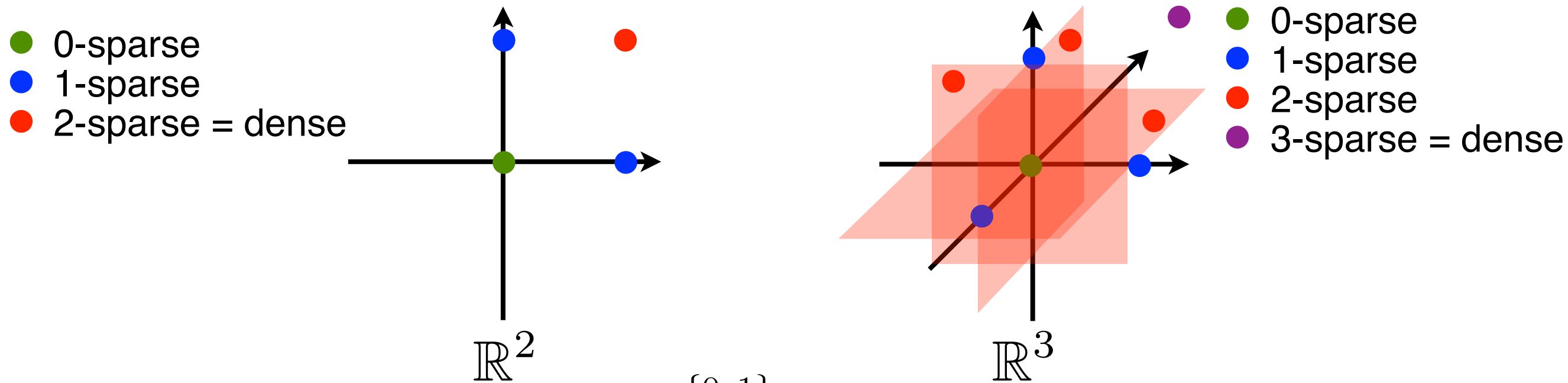


$$\text{supp}(\beta) = \{i = 1, \dots, n : \beta_i \neq 0\}$$

$$\|\beta\|_0 = \#\text{supp}(\beta)$$

(Not a norm : not positively homogenous)

# Blessings in action: Sparsity



Model of  $s$ -sparse vectors : a union of subspaces

$$\mathcal{M}_s = \bigcup_i \{V_i = \text{span}((e_j)_{1 \leq j \leq n}) : \dim(V_i) = s\}.$$

# Models and oracle

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$I^* = \text{supp}(\beta^*)$$

$$\mathcal{I} = \bigcup_{\beta \in \mathbb{R}^p} \text{supp}(\beta)$$

● Known support  $I^*$  :

$$\hat{f} \in X_{I^*} \underset{\beta \in \mathbb{R}^{|I^*|}}{\text{Argmin}} \frac{1}{2\sigma^2} \|y - X_{I^*}\beta\|_2^2 + \frac{n}{2} \log(2\pi\sigma^2)$$

$$\iff \hat{f} = \text{Proj}_{V_{I^*}}(y), \quad V_{I^*} = \text{Span}(X_{I^*}).$$



# Models and oracle

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$I^* = \text{supp}(\beta^*)$$

$$\mathcal{I} = \bigcup_{\beta \in \mathbb{R}^p} \text{supp}(\beta)$$

- Known support  $I^*$  :

$$\hat{f} \in X_{I^*} \underset{\beta \in \mathbb{R}^{|I^*|}}{\text{Argmin}} \frac{1}{2\sigma^2} \|y - X_{I^*}\beta\|_2^2 + \frac{n}{2} \log(2\pi\sigma^2)$$

$$\iff \hat{f} = \text{Proj}_{V_{I^*}}(y), \quad V_{I^*} = \text{Span}(X_{I^*}).$$

- Unknown support  $I^*$  in practice :

- Consider a collection  $\{V_I, I \in \mathcal{I}\}$  of linear subspaces  $\mathbb{R}^n$ , called **models** ;

- Associate to each subspace  $V_I$  the constrained maximum likelihood estimators  $\hat{f}_I = \text{Proj}_{V_I}(y)$  ;

- Estimate  $f^*$  by the **best** estimator among the collection  $\{\hat{f}_I, I \in \mathcal{I}\}$ .

- Meaning of **best** ?

# Models and oracle

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$I^* = \text{supp}(\beta^*)$$

$$\mathcal{I} = \bigcup_{\beta \in \mathbb{R}^p} \text{supp}(\beta)$$

- Risk to measure quality of an estimator :

$$R(\hat{f}) = \mathbb{E} \left[ \|\hat{f} - f^*\|_2^2 \right].$$

- The best estimator in terms of the so-called **oracle** estimator :

$$\hat{f}_{I_b} = \text{Proj}_{V_{I_b}}(y), \quad I_b \in \underset{I \in \mathcal{I}}{\text{Argmin}} R(\hat{f}_I).$$

Bias

Variance

$$R(\hat{f}_I) = \left\| \text{Proj}_{V_I^\perp}(f^*) \right\|_2^2 + \dim(V_I)\sigma^2.$$

- The oracle model  $V_{I_b}$  is that in the collection  $\{V_I : I \in \mathcal{I}\}$  which achieves the best bias-variance trade-off :

- Bias decreases with dimension of  $V_I$ .
- Variance increases with dimension of  $V_I$ .

# Unbiased risk estimator

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon \quad I^* = \text{supp}(\beta^*)$$

$$R(\hat{f}_I) = \underbrace{\| \text{Proj}_{V_I^\perp}(f^*) \|_2^2}_{\text{Bias}} + \underbrace{\dim(V_I)\sigma^2}_{\text{Variance}}. \quad \hat{f}_I = \text{Proj}_{V_I} y$$

- Unfortunately,  $R(\hat{f}_I)$  cannot be computed in practice : depends on  $f^*$  which is unknown to the user.
- Replace  $R(\hat{f}_I)$  by a **good, yet computable**, estimate.
- Observe that

$$\begin{aligned} \mathbb{E} \left[ \|y - \hat{f}_I\|_2^2 \right] &= \mathbb{E} \left[ \| \text{Proj}_{V_I^\perp}(f^* + \varepsilon) \|_2^2 \right] \\ &= \mathbb{E} \left[ \| \text{Proj}_{V_I^\perp} f^* \|_2^2 \right] + \mathbb{E} \left[ \left\langle \text{Proj}_{V_I^\perp} f^*, \varepsilon \right\rangle \right] + \mathbb{E} \left[ \| \text{Proj}_{V_I^\perp} \varepsilon \|_2^2 \right] \\ &= \mathbb{E} \left[ \| \text{Proj}_{V_I^\perp} f^* \|_2^2 \right] + (n - \dim(V_I))\sigma^2 = R(\hat{f}_I) + (n - 2 \dim(V_I))\sigma^2. \end{aligned}$$

- Unbiased risk estimator :

$$\mathbb{E} \left[ \hat{R}(\hat{f}_I) \right] = R(\hat{f}_I), \quad \hat{R}(\hat{f}_I) = \|y - \hat{f}_I\|_2^2 + (2 \dim(V_I) - n)\sigma^2.$$

# Unbiased risk estimator

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon \quad I^* = \text{supp}(\beta^*)$$

$$R(\hat{f}_I) = \underbrace{\| \text{Proj}_{V_I^\perp}(f^*) \|_2^2}_{\text{Bias}} + \underbrace{\dim(V_I)\sigma^2}_{\text{Variance}}. \quad \hat{f}_I = \text{Proj}_{V_I} y$$

- Unfortunately,  $R(\hat{f}_I)$  cannot be computed in practice : depends on  $f^*$  which is unknown to the user.
- Replace  $R(\hat{f}_I)$  by a **good, yet computable**, estimate.
- Observe that

$$\begin{aligned} \mathbb{E} \left[ \|y - \hat{f}_I\|_2^2 \right] &= \mathbb{E} \left[ \| \text{Proj}_{V_I^\perp}(f^* + \varepsilon) \|_2^2 \right] \\ &= \mathbb{E} \left[ \| \text{Proj}_{V_I^\perp} f^* \|_2^2 \right] + \mathbb{E} \left[ \left\langle \text{Proj}_{V_I^\perp} f^*, \varepsilon \right\rangle \right] + \mathbb{E} \left[ \| \text{Proj}_{V_I^\perp} \varepsilon \|_2^2 \right] \\ &= \mathbb{E} \left[ \| \text{Proj}_{V_I^\perp} f^* \|_2^2 \right] + (n - \dim(V_I))\sigma^2 = R(\hat{f}_I) + (n - 2 \dim(V_I))\sigma^2. \end{aligned}$$

- Unbiased risk estimator :

$$\mathbb{E} \left[ \hat{R}(\hat{f}_I) \right] = R(\hat{f}_I), \quad \hat{R}(\hat{f}_I) = \|y - \hat{f}_I\|_2^2 + (2 \dim(V_I) \overset{\text{AIC}}{\cancel{n}})\sigma^2.$$

# Akaike Information Criterion (AIC)

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$I^* = \text{supp}(\beta^*)$$
$$\mathcal{I} = \bigcup_{\beta \in \mathbb{R}^p} \text{supp}(\beta)$$

$$\text{AIC}(\hat{f}_I) = \|\text{Proj}_{V_I^\perp} y\|_2^2 + 2 \dim(V_I) \sigma^2.$$

- Select the model that minimizes AIC :

$$I_{\text{AIC}} \in \underset{I \in \mathcal{I}}{\text{Argmin}} \text{AIC}(\hat{f}_I)$$

- Popular and simple to implement, but ... yields very poor results in high dimension (number of models grows fast with dimension  $p$  as in the sparse case).
- Simple justification assuming  $X = \text{Id}$  :
  - $\text{AIC}(\hat{f}_I) = \|y\|_2^2 + \sum_{i \in I} (2\sigma^2 - |y_i|^2)$ .
  - $I_{\text{AIC}} = \{i : |y_i|^2 > 2\sigma^2\}$  (hard thresholding).
  - If  $f^* = 0$  :  $|I_{\text{AIC}}| \sim \mathcal{B}(p, q)$ ,  $q = \Pr(|Z| > \sqrt{2})$ ,  $Z \sim \mathcal{N}(0, 1)$ .
  - $\mathbb{E}[|I_{\text{AIC}}|] = pq \approx 0.157p$  while  $I_{\text{oracle}} = \emptyset$  !!!



# Bayesian Information Criterion (BIC)

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$I^* = \text{supp}(\beta^*)$$
$$\mathcal{I} = \bigcup_{\beta \in \mathbb{R}^p} \text{supp}(\beta)$$

- Bayesian paradigm,  $I^*$ ,  $f^*$  and  $\varepsilon$  are random :
  - $I^*$  sampled from  $\{\pi_I : I \in \mathcal{I}\}$ ;
  - $f^*$  sampled from  $dP(f|I^*)$  on  $V_{I^*}$ ;
  - Generate  $y$ .
- $P(I|y) \propto \int_{f \in V_I} \pi_I \exp(-\|y - f\|_2^2 / (2\sigma^2)) dP(f|I)$ .
- Posterior likelihood ratio :

$$\log \left( \frac{P(I|y)}{P(I'|y)} \right) \sim_{n \rightarrow +\infty} \frac{\|y - \hat{f}_{I'}\|_2^2 - \|y - \hat{f}_I\|_2^2}{2\sigma^2} + \log(n) \frac{\dim(V_{I'}) - \dim(V_I)}{2}$$
$$+ \log(\pi_I / \pi_{I'}) + O(1).$$

- The most likely model would minimize :

$$\|y - \hat{f}_I\|_2^2 + \sigma^2 \dim(V_I) \log(n) + 2\sigma^2 \log(\pi_I^{-1}).$$

- If the prior  $\pi_I$  is uniform, we get the BIC :

$$I_{\text{BIC}} \in \underset{I \in \mathcal{I}}{\text{Argmin}} \text{BIC}(I) = \|y - \hat{f}_I\|_2^2 + \sigma^2 \dim(V_I) \log(n).$$

# Bayesian Information Criterion (BIC)

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon \quad I^* = \text{supp}(\beta^*)$$
$$\mathcal{I} = \bigcup_{\beta \in \mathbb{R}^p} \text{supp}(\beta)$$

$$I_{\text{BIC}} \in \underset{I \in \mathcal{I}}{\text{Argmin}} \text{BIC}(I) = \|y - \hat{f}_I\|_2^2 + \sigma^2 \dim(V_I) \log(n).$$

- Again, popular and simple poor results in high dimension.
- Simple justification assuming  $X = \text{Id}$  :
  - $I_{\text{BIC}} = \{i : |y_i|^2 > \log(p)\sigma^2\}$  (hard thresholding).
  - If  $f^* = 0$  :  $\mathbb{E} [|I_{\text{BIC}}|] = pq \asymp \sqrt{\frac{2p}{\pi \log(p)}}$ , since

$$q = \Pr(|Z| > \sqrt{\log(p)}) \asymp \sqrt{\frac{2}{\pi \log(p)}} e^{-\log(p)/2} = \sqrt{\frac{2}{\pi p \log(p)}}.$$

$$Z \sim \mathcal{N}(0, 1).$$

- Again  $\mathbb{E} [|I_{\text{BIC}}|]$  grows with  $p$  while  $I_{\text{oracle}} = \emptyset !!!$

# Penalized empirical risk minimization

- How to avoid the selection of a model  $V_I$  with a large dimension ?
- Replace the second penalty term in AIC and BIC by a term taking into account the number of **models per dimension**.
- Associate to the models  $\{V_I : I \in \mathcal{I}\}$  a probability distribution  $\pi = \{\pi_I : I \in \mathcal{I}\}$ .
- Build the model selection criterion ( $K > 1$ , see shortly why)

$$\text{MSC}(I) = \|y - \hat{f}_I\|_2^2 + K\sigma^2 \left( \sqrt{\dim(V_I)} + \sqrt{2 \log(\pi_I^{-1})} \right)^2$$

$$I_{\text{MSC}} \in \underset{I \in \mathcal{I}}{\text{Argmin}} \text{MSC}(I).$$

- Such a choice of penalty seems cryptic but ...
- it ensures that  $R(\hat{f}_{I_{\text{MSC}}})$  is close to  $R(\hat{f}_{I_{\text{oracle}}})$  :
- $\pi$  is chosen to penalize overly high-dimensional models.

# Penalized empirical risk minimization

- Put the same mass on all models with the same sparsity level.
- Choice 1 :  $\pi_I = (1 + 1/p)^{-p} p^{-|I|}$ , and thus

$$\log(\pi_I^{-1}) \leq 1 + |I| \log(p).$$

- Choice 2 :  $\pi_I = (e - 1)/(e - e^{-p})(C_p^{|I|})^{-1} e^{-|I|}$ , and thus

$$\log(\pi_I^{-1}) \leq \log(e/(e - 1)) + 2|I|(2 + \log(p/|I|)).$$

# Penalized empirical risk minimization

- Put the same mass on all models with the same sparsity level.
- Choice 1 :  $\pi_I = (1 + 1/p)^{-p} p^{-|I|}$ , and thus

$$\log(\pi_I^{-1}) \leq 1 + |I| \log(p).$$

- Choice 2 :  $\pi_I = (e - 1)/(e - e^{-p})(C_p^{|I|})^{-1} e^{-|I|}$ , and thus

$$\log(\pi_I^{-1}) \leq \log(e/(e - 1)) + 2|I|(2 + \log(p/|I|)).$$

- Conclusion :  $\log(\pi_I^{-1})$  comparable to  $\dim(V_I)$  (up to the log factor).
- The  $\log(p)$  factor reflects the number of models per sparsity  $|I|$  :  $\log(C_p^{|I|}) \lesssim |I| \log(p/|I|) \leq |I| \log(p)$ .
- It can be shown that this factor is unavoidable.
- For  $X = \text{Id}$ , it can be shown that

$$\mathbb{E} [ |I_{\text{MSC}}| ] \asymp \frac{p^{1-K}}{\sqrt{\pi K \log(p)}} \rightarrow 0 \text{ if } K > 1.$$

# Penalized empirical risk minimization

- Put the same mass on all models with the same sparsity level.
- Choice 1 :  $\pi_I = (1 + 1/p)^{-p} p^{-|I|}$ , and thus

$$\log(\pi_I^{-1}) \leq 1 + |I| \log(p).$$

- Choice 2 :  $\pi_I = (e - 1)/(e - e^{-p})(C_p^{|I|})^{-1} e^{-|I|}$ , and thus

$$\log(\pi_I^{-1}) \leq \log(e/(e - 1)) + 2|I|(2 + \log(p/|I|)).$$

- Conclusion :  $\log(\pi_I^{-1})$  comparable to  $\dim(V_I)$  (up to the log factor).
- The  $\log(p)$  factor reflects the number of models per sparsity  $|I|$  :  $\log(C_p^{|I|}) \lesssim |I| \log(p/|I|) \leq |I| \log(p)$ .
- It can be shown that this factor is unavoidable.
- For  $X = \text{Id}$ , it can be shown that

$$\mathbb{E} [||I_{\text{MSC}}||] \asymp \frac{p^{1-K}}{\sqrt{\pi K \log(p)}} \rightarrow 0 \text{ if } K > 1.$$



---

# ***Guarantees***

# Oracle inequality

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon \quad \mathcal{I} = \bigcup_{\beta \in \mathbb{R}^p} \text{supp}(\beta)$$

$$I_{\text{MSC}} \in \underset{I \in \mathcal{I}}{\text{Argmin}} \left\| y - \hat{f}_I \right\|_2^2 + K \sigma^2 \left( \sqrt{\dim(V_I)} + \sqrt{2 \log(\pi_I^{-1})} \right)^2$$



# Oracle inequality

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon \quad \mathcal{I} = \bigcup_{\beta \in \mathbb{R}^p} \text{supp}(\beta)$$

$$I_{\text{MSC}} \in \underset{I \in \mathcal{I}}{\text{Argmin}} \left\| y - \hat{f}_I \right\|_2^2 + K\sigma^2 \left( \sqrt{\dim(V_I)} + \sqrt{2 \log(\pi_I^{-1})} \right)^2$$

**Theorem** *There exists a constant  $C > 1$  depending only on  $K > 1$ , such that*

$$\mathbb{E} \left[ \left\| \hat{f}_{I_{\text{MSC}}} - f^* \right\|_2^2 \right] \leq C \min_{I \in \mathcal{I}} \left[ \overset{\text{Risk}}{\mathbb{E} \left[ \left\| \hat{f}_I - f^* \right\|_2^2 \right]} + \overset{\text{Complexity}}{\sigma^2 \log(\pi_I^{-1})} + \sigma^2 \right].$$

# Oracle inequality

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon \quad \mathcal{I} = \bigcup_{\beta \in \mathbb{R}^p} \text{supp}(\beta)$$

$$I_{\text{MSC}} \in \underset{I \in \mathcal{I}}{\text{Argmin}} \left\| y - \hat{f}_I \right\|_2^2 + K\sigma^2 \left( \sqrt{\dim(V_I)} + \sqrt{2 \log(\pi_I^{-1})} \right)^2$$

**Theorem** *There exists a constant  $C > 1$  depending only on  $K > 1$ , such that*

$$\mathbb{E} \left[ \left\| \hat{f}_{I_{\text{MSC}}} - f^* \right\|_2^2 \right] \leq C \min_{I \in \mathcal{I}} \left[ \overset{\text{Risk}}{\mathbb{E} \left[ \left\| \hat{f}_I - f^* \right\|_2^2 \right]} + \overset{\text{Complexity}}{\sigma^2 \log(\pi_I^{-1})} + \sigma^2 \right].$$

● **Choice 1** :  $\pi_I = (1 + 1/p)^{-p} p^{-|I|}$ , and thus

$$\mathbb{E} \left[ \left\| \hat{f}_{I_{\text{MSC}}} - f^* \right\|_2^2 \right] \leq C \min_{I \in \mathcal{I}} \left[ \mathbb{E} \left[ \left\| \hat{f}_I - f^* \right\|_2^2 \right] + \overset{\text{Complexity}}{3\sigma^2 |I| \log(p)} \right].$$

● **Choice 2** :  $\pi_I = (e - 1)/(e - e^{-p})(C_p^{|I|})^{-1} e^{-|I|}$ , and thus

$$\mathbb{E} \left[ \left\| \hat{f}_{I_{\text{MSC}}} - f^* \right\|_2^2 \right] \leq C \min_{I \in \mathcal{I}} \left[ \mathbb{E} \left[ \left\| \hat{f}_I - f^* \right\|_2^2 \right] + \overset{\text{Complexity}}{3\sigma^2 |I| \log(p/|I|)} \right].$$

● **Conclusion** : MSC achieves the best tradeoff between the risk and the complexity of the model.

# Optimality: minimax risk

$$y = \underbrace{X\beta^*}_{f^* \in \mathcal{M}_s} + \varepsilon \quad \mathcal{M}_s \stackrel{\text{def}}{=} \bigcup_{I \in \mathcal{I}: |I|=s} \underbrace{\text{Span}(X_I)}_{V_I}, \quad s \leq p/2$$

$$\begin{aligned} \sup_{f^* \in \mathcal{M}_s} \mathbb{E} \left[ \|\hat{f}_{I_{\text{MSC}}} - f^*\|_2^2 \right] &\leq C \sup_{f^* \in \mathcal{M}_s} \min_{I \in \mathcal{I}} \left[ \mathbb{E} \left[ \|\hat{f}_I - f^*\|_2^2 \right] + 3\sigma^2 |I| \log(p/|I|) \right] \\ &\leq C \sup_{f^* \in \mathcal{M}_s} \min_{I \in \mathcal{I}: |I|=s} \left[ \underbrace{\|\text{Proj}_{V_I^\perp} f^*\|_2^2}_{=0} + 4\sigma^2 s \log(p/s) \right] = C' \sigma^2 s \log(p/s). \end{aligned}$$

# Optimality: minimax risk

$$y = \underbrace{X\beta^*}_{f^* \in \mathcal{M}_s} + \varepsilon \quad \mathcal{M}_s \stackrel{\text{def}}{=} \bigcup_{I \in \mathcal{J}: |I|=s} \underbrace{\text{Span}(X_I)}_{V_I}, \quad s \leq p/2$$

$$\begin{aligned} \sup_{f^* \in \mathcal{M}_s} \mathbb{E} \left[ \|\hat{f}_{I_{\text{MSC}}} - f^*\|_2^2 \right] &\leq C \sup_{f^* \in \mathcal{M}_s} \min_{I \in \mathcal{J}} \left[ \mathbb{E} \left[ \|\hat{f}_I - f^*\|_2^2 \right] + 3\sigma^2 |I| \log(p/|I|) \right] \\ &\leq C \sup_{f^* \in \mathcal{M}_s} \min_{I \in \mathcal{J}: |I|=s} \left[ \underbrace{\|\text{Proj}_{V_I^\perp} f^*\|_2^2}_{=0} + 4\sigma^2 s \log(p/s) \right] = C' \sigma^2 s \log(p/s). \end{aligned}$$

$$\delta_{\min} = \inf_{\beta: \|\beta\|_0=2s} \frac{\|X\beta\|_2}{\|\beta\|_2} \leq \sup_{\beta: \|\beta\|_0=2s} \frac{\|X\beta\|_2}{\|\beta\|_2} = \delta_{\max}.$$

**Theorem** For any  $s \leq p/5$ , we have

$$\inf_{\hat{f}} \sup_{f^* \in \mathcal{M}_s} \mathbb{E} \left[ \|\hat{f} - f^*\|_2^2 \right] \geq \frac{e}{8(2e+1)^2} \left( \frac{\delta_{\min}}{\delta_{\max}} \right)^2 \sigma^2 s \log(p/(5s)).$$

# Optimality: minimax risk

$$y = \underbrace{X\beta^*}_{f^* \in \mathcal{M}_s} + \varepsilon \quad \mathcal{M}_s \stackrel{\text{def}}{=} \bigcup_{I \in \mathcal{I}: |I|=s} \underbrace{\text{Span}(X_I)}_{V_I}, \quad s \leq p/2$$

$$\begin{aligned} \sup_{f^* \in \mathcal{M}_s} \mathbb{E} \left[ \|\hat{f}_{I_{\text{MSC}}} - f^*\|_2^2 \right] &\leq C \sup_{f^* \in \mathcal{M}_s} \min_{I \in \mathcal{I}} \left[ \mathbb{E} \left[ \|\hat{f}_I - f^*\|_2^2 \right] + 3\sigma^2 |I| \log(p/|I|) \right] \\ &\leq C \sup_{f^* \in \mathcal{M}_s} \min_{I \in \mathcal{I}: |I|=s} \left[ \underbrace{\|\text{Proj}_{V_I^\perp} f^*\|_2^2}_{=0} + 4\sigma^2 s \log(p/s) \right] = C' \sigma^2 s \log(p/s). \end{aligned}$$

$$\delta_{\min} = \inf_{\beta: \|\beta\|_0=2s} \frac{\|X\beta\|_2}{\|\beta\|_2} \leq \sup_{\beta: \|\beta\|_0=2s} \frac{\|X\beta\|_2}{\|\beta\|_2} = \delta_{\max}.$$

**Theorem** For any  $s \leq p/5$ , we have

$$\inf_{\hat{f}} \sup_{f^* \in \mathcal{M}_s} \mathbb{E} \left[ \|\hat{f} - f^*\|_2^2 \right] \geq \frac{e}{8(2e+1)^2} \left( \frac{\delta_{\min}}{\delta_{\max}} \right)^2 \sigma^2 s \log(p/(5s)).$$

**Upper and lower bound are of the same order  $\sim s \log(n/s) \Rightarrow \hat{f}_{I_{\text{MSC}}}$  is optimal.**

# Optimality: minimax risk

$$y = \underbrace{X\beta^*}_{f^* \in \mathcal{M}_s} + \varepsilon \quad \mathcal{M}_s \stackrel{\text{def}}{=} \bigcup_{I \in \mathcal{I}: |I|=s} \underbrace{\text{Span}(X_I)}_{V_I}, \quad s \leq p/2$$

$$\begin{aligned} \sup_{f^* \in \mathcal{M}_s} \mathbb{E} \left[ \|\hat{f}_{I_{\text{MSC}}} - f^*\|_2^2 \right] &\leq C \sup_{f^* \in \mathcal{M}_s} \min_{I \in \mathcal{I}} \left[ \mathbb{E} \left[ \|\hat{f}_I - f^*\|_2^2 \right] + 3\sigma^2 |I| \log(p/|I|) \right] \\ &\leq C \sup_{f^* \in \mathcal{M}_s} \min_{I \in \mathcal{I}: |I|=s} \left[ \underbrace{\|\text{Proj}_{V_I^\perp} f^*\|_2^2}_{=0} + 4\sigma^2 s \log(p/s) \right] = C' \sigma^2 s \log(p/s). \end{aligned}$$

$$\delta_{\min} = \inf_{\beta: \|\beta\|_0=2s} \frac{\|X\beta\|_2}{\|\beta\|_2} \leq \sup_{\beta: \|\beta\|_0=2s} \frac{\|X\beta\|_2}{\|\beta\|_2} = \delta_{\max}.$$

**Theorem** For any  $s \leq p/5$ , we have

$$\inf_{\hat{f}} \sup_{f^* \in \mathcal{M}_s} \mathbb{E} \left[ \|\hat{f} - f^*\|_2^2 \right] \geq \frac{e}{8(2e+1)^2} \left( \frac{\delta_{\min}}{\delta_{\max}} \right)^2 \sigma^2 s \log(p/(5s)).$$

**Upper and lower bound are of the same order  $\sim s \log(n/s) \Rightarrow \hat{f}_{I_{\text{MSC}}}$  is optimal.**

**Rmainder captures a phase transition phenomenon between good and poor estimation, i.e.  $n \gtrsim s \log(p)$ .**

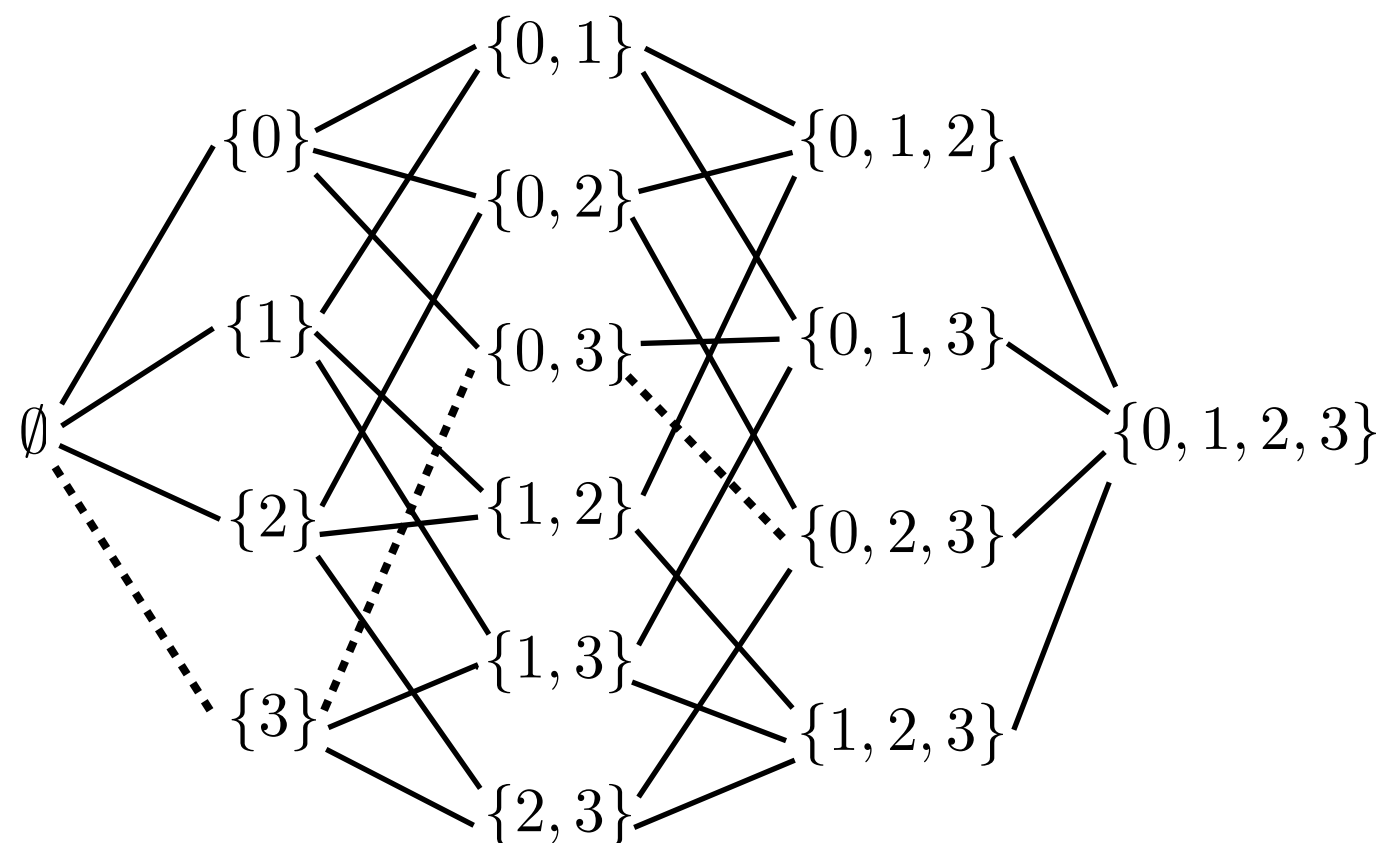
# Computational issues

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$\mathcal{I} = \bigcup_{\beta \in \mathbb{R}^p} \text{supp}(\beta)$$
$$|\mathcal{I}| = 2^p$$

$$I_{\text{MSC}} \in \underset{I \in \mathcal{I}}{\text{Argmin}} \left\| y - \hat{f}_I \right\|_2^2 + K\sigma^2 \left( \sqrt{\dim(V_I)} + \sqrt{2 \log(\pi_I^{-1})} \right)^2$$

**This cannot be implemented in practice except for the orthogonal case (requires exhaustive search) : NP-hard problem.**



# Computational issues

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon \quad \mathcal{I} = \bigcup_{\beta \in \mathbb{R}^p} \text{supp}(\beta)$$
$$|\mathcal{I}| = 2^p$$

$$I_{\text{MSC}} \in \underset{I \in \mathcal{I}}{\text{Argmin}} \left\| y - \hat{f}_I \right\|_2^2 + K\sigma^2 \left( \sqrt{\dim(V_I)} + \sqrt{2 \log(\pi_I^{-1})} \right)^2$$

**This cannot be implemented in practice except for the orthogonal case (requires exhaustive search) : NP-hard problem.**

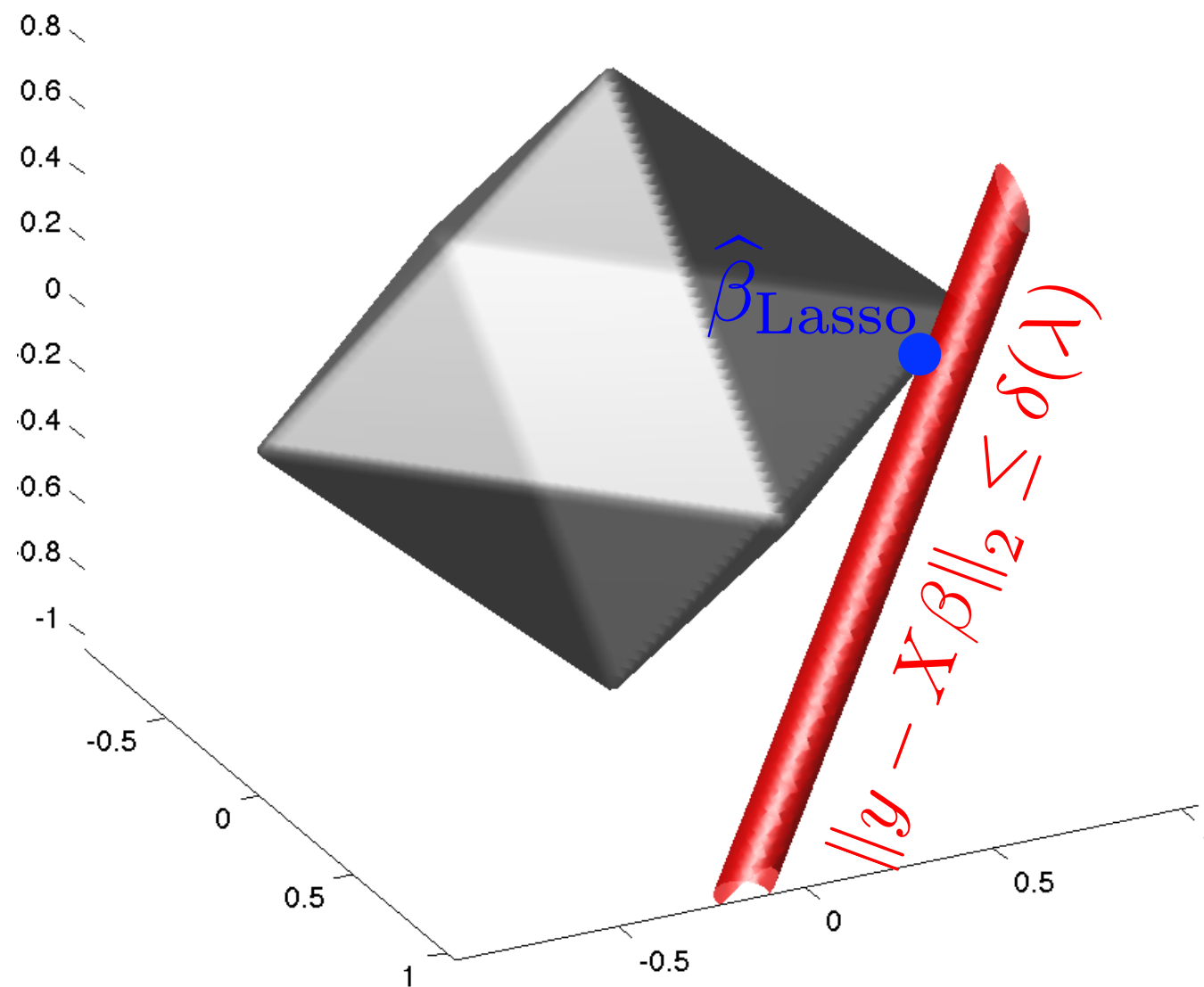
**Can we have a polynomial-time procedure to do so with the same performance guarantees ?**



# Convex relaxation: Lasso

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$\hat{f}_{\text{Lasso}} = X \underset{\beta \in \mathbb{R}^p}{\text{Argmin}} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$



# Lasso oracle inequality

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$\hat{f}_{\text{Lasso}} = X \underset{\beta \in \mathbb{R}^p}{\text{Argmin}} \left\| y - X\beta \right\|_2^2 + \lambda \|\beta\|_1$$

$$\Upsilon(I, c) \stackrel{\text{def}}{=} \inf_{\{\alpha \in \mathbb{R}^p : \|\alpha_{I^c}\|_1 < c\|\alpha_I\|_1\}} \frac{\sqrt{|I|} \|X\alpha\|_2}{\|\alpha_I\|_1} \quad \text{A measure of restricted ill-conditioning}$$

# Lasso oracle inequality

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$\hat{f}_{\text{Lasso}} = X \underset{\beta \in \mathbb{R}^p}{\text{Argmin}} \left\| y - X\beta \right\|_2^2 + \lambda \|\beta\|_1$$

$$\Upsilon(I, c) \stackrel{\text{def}}{=} \inf_{\{\alpha \in \mathbb{R}^p : \|\alpha_{I^c}\|_1 < c \|\alpha_I\|_1\}} \frac{\sqrt{|I|} \|X\alpha\|_2}{\|\alpha_I\|_1} \quad \text{A measure of restricted ill-conditioning}$$

**Theorem** Assume that the columns of  $X$  have unit-norm. Then,  $\forall \delta > 1$ , the Lasso with

$$\lambda = 2\sigma \sqrt{2\delta \log(2p)},$$

obeys with probability at least  $1 - 2(2p)^{1-\delta}$

$$\|\hat{f}_{\text{Lasso}} - f^*\|_2^2 \leq \min_{\substack{I \in \mathcal{I} \\ \text{supp}(\beta) = I}} \left[ \left\| X\beta - f^* \right\|_2^2 + \frac{\text{Complexity}}{\Upsilon(I, 3)^2} |I| \log(2p) \right].$$

# Lasso oracle inequality

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$\hat{f}_{\text{Lasso}} = X \underset{\beta \in \mathbb{R}^p}{\text{Argmin}} \left\| y - X\beta \right\|_2^2 + \lambda \left\| \beta \right\|_1$$

$$\Upsilon(I, c) \stackrel{\text{def}}{=} \inf_{\{\alpha \in \mathbb{R}^p : \|\alpha_{I^c}\|_1 < c \|\alpha_I\|_1\}} \frac{\sqrt{|I|} \|X\alpha\|_2}{\|\alpha_I\|_1} \quad \text{A measure of restricted ill-conditioning}$$

**Theorem** Assume that the columns of  $X$  have unit-norm. Then,  $\forall \delta > 1$ , the Lasso with

$$\lambda = 2\sigma \sqrt{2\delta \log(2p)},$$

obeys with probability at least  $1 - 2(2p)^{1-\delta}$

$$\left\| \hat{f}_{\text{Lasso}} - f^* \right\|_2^2 \leq \min_{\substack{I \in \mathcal{J} \\ \text{supp}(\beta) = I}} \left[ \left\| X\beta - f^* \right\|_2^2 + \frac{\text{Complexity}}{\Upsilon(I, 3)^2} |I| \log(2p) \right].$$

$$\mathbb{E} \left[ \left\| \hat{f}_{I_{\text{MSC}}} - f^* \right\|_2^2 \right] \leq C \min_{I \in \mathcal{J}} \left[ \mathbb{E} \left[ \left\| \hat{f}_I - f^* \right\|_2^2 \right] + 3\sigma^2 |I| \log(p) \right].$$

**Remainder term of the same order as MSC while Lasso is implementable.**

**The OI is sharp for Lasso ( $C = 1$ ).**

# Implementation of Lasso

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon \quad \hat{f}_{\text{Lasso}} = X \underset{\beta \in \mathbb{R}^p}{\text{Argmin}} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

***See MC on non-smooth optimization.***

# Aggregation by exponential weighting

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

- Idea : rather than selecting one model, take the best of each model by averaging.
- The Exponential Weighted Aggregation (EWA) estimator is

$$F(\beta) \stackrel{\text{def}}{=} \|y - X\beta\|_2^2 + \lambda\|\beta\|_1 \quad \hat{\beta}_{\text{EWA}} = \int_{\mathbb{R}^p} \beta \mu(\beta) d\beta$$

$$\mu(\beta) = \frac{\exp(-F(\beta)/\tau)}{\int_{\mathbb{R}^p} \exp(-F(\alpha)/\tau) d\alpha},$$

# Aggregation by exponential weighting

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

- Idea : rather than selecting one model, take the best of each model by averaging.
- The Exponential Weighted Aggregation (EWA) estimator is

$$F(\beta) \stackrel{\text{def}}{=} \|y - X\beta\|_2^2 + \lambda\|\beta\|_1 \quad \hat{\beta}_{\text{EWA}} = \int_{\mathbb{R}^p} \beta \mu(\beta) d\beta$$

$$\mu(\beta) = \frac{\exp(-F(\beta)/\tau)}{\int_{\mathbb{R}^p} \exp(-F(\alpha)/\tau) d\alpha},$$

Involves solving an **integration** problem.

# EWA oracle inequality

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$\hat{\beta}_{\text{EWA}} = \int_{\mathbb{R}^p} \beta \mu(\beta) d\beta$$

$$\Upsilon(I, c) \stackrel{\text{def}}{=} \inf_{\{\alpha \in \mathbb{R}^p : \|\alpha_{I^c}\|_1 < c \|\alpha_I\|_1\}} \frac{\sqrt{|I|} \|X\alpha\|_2}{\|\alpha_I\|_1} \quad \text{A measure of restricted ill-conditioning}$$



# EWA oracle inequality

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$\hat{\beta}_{\text{EWA}} = \int_{\mathbb{R}^p} \beta \mu(\beta) d\beta$$

$$\Upsilon(I, c) \stackrel{\text{def}}{=} \inf_{\{\alpha \in \mathbb{R}^p : \|\alpha_{I^c}\|_1 < c \|\alpha_I\|_1\}} \frac{\sqrt{|I|} \|X\alpha\|_2}{\|\alpha_I\|_1} \quad \text{A measure of restricted ill-conditioning}$$

**Theorem** Assume that the columns of  $X$  have unit-norm. Then,  $\forall \delta > 1$ , and for some absolute constant  $C' > 0$ , the EWA with

$$\lambda = 2\sigma \sqrt{2\delta \log(2p)} \quad \text{and} \quad \tau = O(1/p)$$

obeys with probability at least  $1 - 2(2p)^{1-\delta}$

$$\|\hat{f}_{\text{Lasso}} - f^*\|_2^2 \leq \min_{\substack{I \in \mathcal{I} \\ \text{supp}(\beta) = I}} \left[ \|X\beta - f^*\|_2^2 + \frac{\text{Complexity}}{\Upsilon(I, 3)^2} |I| \log(2p) \right],$$

# EWA oracle inequality

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$\hat{\beta}_{\text{EWA}} = \int_{\mathbb{R}^p} \beta \mu(\beta) d\beta$$

$$\Upsilon(I, c) \stackrel{\text{def}}{=} \inf_{\{\alpha \in \mathbb{R}^p : \|\alpha_{I^c}\|_1 < c \|\alpha_I\|_1\}} \frac{\sqrt{|I|} \|X\alpha\|_2}{\|\alpha_I\|_1} \quad \text{A measure of restricted ill-conditioning}$$

**Theorem** Assume that the columns of  $X$  have unit-norm. Then,  $\forall \delta > 1$ , and for some absolute constant  $C' > 0$ , the EWA with

$$\lambda = 2\sigma \sqrt{2\delta \log(2p)} \quad \text{and} \quad \tau = O(1/p)$$

obeys with probability at least  $1 - 2(2p)^{1-\delta}$

$$\begin{aligned} \|\hat{f}_{\text{Lasso}} - f^*\|_2^2 &\leq \min_{\substack{I \in \mathcal{J} \\ \text{supp}(\beta) = I}} \left[ \|X\beta - f^*\|_2^2 + \frac{\text{Complexity}}{\Upsilon(I, 3)^2} |I| \log(2p) \right], \\ \mathbb{E} \left[ \|\hat{f}_{I_{\text{MSC}}} - f^*\|_2^2 \right] &\leq C \min_{I \in \mathcal{J}} \left[ \mathbb{E} \left[ \|\hat{f}_I - f^*\|_2^2 \right] + 3\sigma^2 |I| \log(p) \right]. \end{aligned}$$

**Remainder term of the same order as MSC and Lasso.**

**The OI is sharp for EWA as well.**

# Implementation of EWA

$$y = \underbrace{X\beta^*}_{f^*} + \varepsilon$$

$$\hat{\beta}_{\text{EWA}} = \int_{\mathbb{R}^p} \beta \mu(\beta) d\beta$$

- An integration problem in high dimension.
- Very challenging.
- MC sampling through e.g. SDE's.
- Maybe another MC.

# Take-away messages

- Model selection benefits from blessings of dimensionality.
- Several model selection criteria in the literature.
- AIC and BIC are not good when the number of models is very large.
- MSC is much better and enjoys nice guarantees, but is NP-hard.
- Lasso is implementable (convex program) while enjoying the same guarantees.
- So does EWA but necessitates even more sophisticated sampling algorithms.
- We gave principles for estimation but other tasks as well: classification, machine learning.

---

<https://fadili.users.greyc.fr/>

**Thanks**  
**Any questions ?**