

# Linear regression with random projections (Gaussian weak features model)

Belkin, Hsu, Xu (2019) arXiv: 1903.07571

Breiman, Friedman (1983) JASA vol 78, 361

Lafor, Thomas (2021) arXiv: 2403.10459

Consider a data set  $S = \{(\vec{x}^k, y^k)\}_{k=1}^n$   
with each pair  $(\vec{x}^k, y^k) \in \mathbb{R}^{d+1}$  sampled  
i.i.d from a distribution  $\mathcal{D}(\vec{x}, y)$ . Assume:

$$\mathcal{D}(\vec{x}) = \mathcal{N}(\vec{x} | \vec{0}, \mathbb{I}_d)$$

$\mathcal{D}(y | \vec{x})$ : modeled by the linear function

$$y = \vec{\beta}^T \vec{x} + \epsilon; \quad \epsilon \sim \mathcal{N}(\epsilon | 0, 1)$$

"ground-truth vector":  $\vec{\beta} \in \mathbb{R}^d$

let:

$$X \equiv \begin{bmatrix} \vec{x}^1 & \dots & \vec{x}^n \end{bmatrix}^T \in \mathbb{R}^{n \times d}$$
$$\vec{y} \equiv \begin{bmatrix} y^1 & y^2 & \dots & y^n \end{bmatrix}^T \in \mathbb{R}^n$$

# Least-squares estimation

In usual least-squares estimation, one searches a minimizer  $\hat{\beta} \in \mathbb{R}^d$  for the empirical loss quadratic loss:

$$\begin{aligned}\hat{\beta} &= \operatorname{argmin}_{\beta \in \mathbb{R}^d} \mathcal{L}_S(S; \beta) \\ &= \operatorname{argmin}_{\beta \in \mathbb{R}^d} \underbrace{\sum_{k=1}^n (y^k - \beta^T x^k)^2}_{\text{squared loss function}}\end{aligned}$$

$$= \operatorname{argmin}_{\beta \in \mathbb{R}^d} (\vec{y} - X\beta)^T (\vec{y} - X\beta)$$

$$= \operatorname{argmin}_{\beta \in \mathbb{R}^d} \|\vec{y} - X\beta\|^2$$

In the exercises we also discuss a particular kind of regularized least squares.

## linear regression with (random) projections

We will fit a linear model to the data using only a subset:

$A \subseteq [d] \equiv \{1, \dots, d\}$  of  $p \equiv |A|$  variables

For any  $\vec{v} \in \mathbb{R}^d$  we use

$$\vec{v}_A \equiv [v_f : f \in A]^T$$

to denote its  $|A|$ -dimensional subvector of entries from  $A$ . Also denote

$$\vec{X}_A \equiv [\vec{x}_A^1 \mid \dots \mid \vec{x}_A^n]^T \in \mathbb{R}^{n \times |A|}$$

For  $A \subseteq [d]$ , its complement is denoted by  $A^c \equiv [d] \setminus A$ .

The regression coefficients  $\hat{\beta} \in \mathbb{R}^d$  are fitted with

$$\hat{\beta}_A \equiv X_A^T \hat{y}$$
$$\hat{\beta}_{A^c} \equiv \mathbf{0}$$

• Solution for the least-squares problem

$$X_A^T X_A \hat{\beta}_A = X_A^T \hat{y}$$

for  $\hat{\beta}_A$

•  $\hat{\beta}_A^c$  forced to all-zeros

$\dagger$  : denotes the Moore-Penrose inverse

Definition: Moore-Penrose inverse

For  $M \in \mathbb{R}^{l \times s}$ , a Moore-Penrose or pseudo inverse is defined as the matrix  $M^\dagger \in \mathbb{R}^{s \times l}$  satisfying all the three criteria:



1)  $MM^T$  need not be the identity matrix, but it maps all the column vectors of  $M$  to themselves:

$$MM^T M = M$$

2)  $M^T$  acts like a weak inverse:

$$M^T M M^T = M^T$$

3)  $MM^T$  and  $M^T M$  are symmetric:

$$(MM^T)^T = MM^T ; (M^T M)^T = M^T M$$

Note that  $MM^T$  and  $M^T M$  are orthogonal projection operators, as follows from

$$(MM^T)^2 = MM^T ; (M^T M)^2 = M^T M$$

## Test risk / generalization error

Once the estimator  $\hat{\beta} \in \mathbb{R}^d$  is computed, its quality over a new sample pair:

$$\left( \vec{x}^{\text{new}}, y^{\text{new}} \right) \sim \mathcal{D} \leftarrow \begin{array}{l} \text{same distribution} \\ \text{that generated} \\ \text{the training} \\ \text{data} \end{array}$$

can be measured by:

$$l(\vec{x}^{\text{new}}, y^{\text{new}}; \hat{\beta})$$

prediction error

Conditional mean prediction risk

$$\mathcal{E}(X, \vec{y}; \hat{\beta}) \equiv \mathbb{E}_{\vec{x}^{\text{new}}, y^{\text{new}} | X, \vec{y}} \left[ l(\vec{x}^{\text{new}}, y^{\text{new}}; \hat{\beta}) \right]$$

The test risk is then defined by:

$$\mathcal{R}(\hat{\beta}) \equiv \mathbb{E}_{X, \vec{y}} [E(X, \vec{y}; \hat{\beta})]$$

↳ expectation of conditional prediction error over the data distribution  
a.k.a. population risk

## Theorem

Assume  $\vec{x} \sim \mathcal{N}(\vec{x} | \vec{0}, \mathbb{I}_d)$ ;  $\epsilon \sim \mathcal{N}(\epsilon | 0, 1)$  independent of  $\vec{x}$  and  $y = \vec{\beta}^T \vec{x} + \epsilon$  for some  $\vec{\beta} \in \mathbb{R}^d$  and  $\mu > 0$ . Pick any  $p \in \{0, \dots, d\}$  and  $A \subseteq [d]$  with  $|A| = p$ .

For the squared loss:

$$l(\vec{x}, y; \hat{\beta}) = (y - \hat{\beta}^T \vec{x})^2$$

with  $\hat{\beta}_A = X_A^T \vec{y}$  and  $\hat{\beta}_{A^c} = \vec{0}$ , the test

risk  $\mathcal{R}_A(\vec{\beta})$  of  $\hat{\vec{\beta}}$  for a given  $A$  is:

$$\mathcal{R}_A(\vec{\beta}) =$$

$$\begin{cases} \left( \|\vec{\beta}_{A^c}\|^2 + \eta^2 \right) \left( 1 + \frac{p}{n-p-1} \right); & \text{if } p \leq n-2 \\ + \infty; & \text{if } n-1 \leq p \leq n+1 \\ \|\vec{\beta}_A\|^2 \left( 1 - \frac{n}{p} \right) + \left( \|\vec{\beta}_{A^c}\|^2 + \eta^2 \right) \left( 1 + \frac{n}{p-n-1} \right); & \text{if } p \geq n+2 \end{cases}$$

### Corollary

Let  $A$  be a uniformly random subset of  $[d]$  of cardinality  $p$ . In the setting of the theorem above, we have:

$$\mathcal{R}(\vec{\beta}) \equiv \mathbb{E}_A \left[ \mathcal{R}_A(\vec{\beta}) \right] =$$

$$\left\{ \begin{aligned} & \left( \left( 1 - \frac{p}{d} \right) \|\vec{\beta}\|^2 + \eta^2 \right) \left( 1 + \frac{p}{n-p-1} \right) ; \\ & \quad \text{if } p \leq n-z \\ & \|\vec{\beta}\|^2 \left[ 1 - \frac{n}{d} \left( z - \frac{d-n-1}{p-n-1} \right) \right] \\ & + \eta^2 \left( 1 + \frac{n}{p-n-1} \right) ; \quad \text{if } p \geq n+z \end{aligned} \right.$$

### Proof of the corollary

Since  $A$  is a uniformly random subset of  $[d]$  of cardinality  $p$ :

$$E_A [\|\vec{\beta}_A\|^2] = \frac{p}{d} \|\vec{\beta}\|^2$$

$$E_{A^c} [\|\vec{\beta}_{A^c}\|^2] = \left( 1 - \frac{p}{d} \right) \|\vec{\beta}\|^2$$

Plugging into Theorem 1 completes the proof.

# Sketching the plot for the corollary

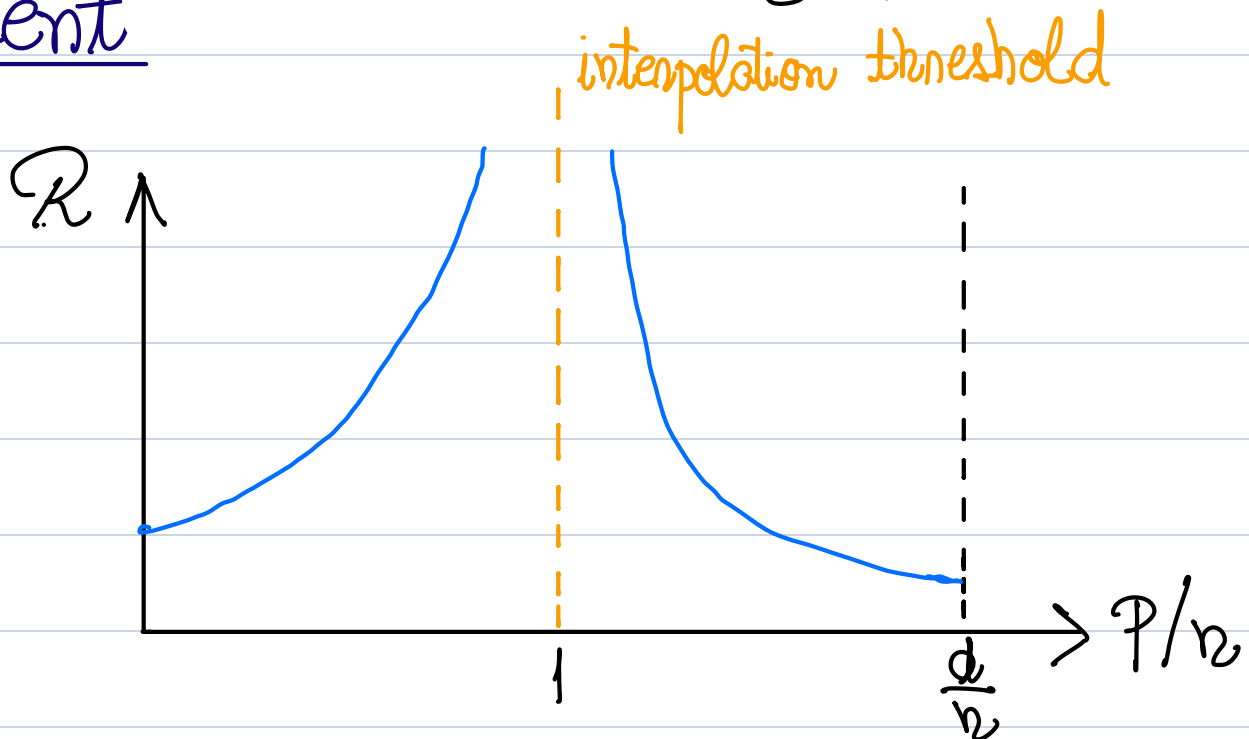
Assume  $d > n+1$ :

a) The risk first increases with  $p$  up to the "interpolation threshold"  $\underline{p=n}$ , after which decreases with  $p$ .

b) If

$$\frac{\|\vec{\beta}\|^2}{y^2} > \frac{d}{d-n-1},$$

the smallest test risk is achieved at  $p=d$ .  
It is smaller than any  $p \leq n$ : double descent



# Proof of the theorem

Let us consider the conditional mean squared risk:

$$\begin{aligned} \mathcal{E}(X, \vec{y}; \hat{\vec{\beta}}) &= \mathbb{E}_{\vec{x}^{\text{new}}, y^{\text{new}}} \left[ (y^{\text{new}} - \hat{\vec{\beta}}^T \vec{x}^{\text{new}})^2 \right] = \\ &= \mathbb{E}_{\vec{x}^{\text{new}}, y^{\text{new}}} \left[ \left( (\vec{\beta} - \hat{\vec{\beta}})^T \vec{x}^{\text{new}} + y \epsilon^{\text{new}} \right)^2 \right] = \\ &= (\vec{\beta} - \hat{\vec{\beta}})^T \underbrace{\mathbb{E}_{\vec{x}^{\text{new}}} \left[ \vec{x}^{\text{new}} \vec{x}^{\text{new}T} \right]}_{\mathbb{I}_d} (\vec{\beta} - \hat{\vec{\beta}}) \\ &\quad + \underbrace{\mu^2 \mathbb{E}_{\epsilon^{\text{new}}} \left[ (\epsilon^{\text{new}})^2 \right]}_{=1} + \underbrace{2\mu (\vec{\beta} - \hat{\vec{\beta}})^T \mathbb{E}_{\vec{x}^{\text{new}}} \left[ \vec{x}^{\text{new}} \right]}_{=0} \underbrace{\mathbb{E}_{\epsilon^{\text{new}}} \left[ \epsilon^{\text{new}} \right]}_{=0} \\ &= \sum_{f \in A} (\beta_f - \hat{\beta}_f)^2 + \sum_{e \in A^c} (\beta_e - \hat{\beta}_e)^2 + \mu^2 \end{aligned}$$

Since  $\hat{\vec{\beta}}_{A^c} = \vec{0}$  we have:

$$\mathcal{E}(X, \vec{y}; \hat{\vec{\beta}}) = \mu^2 + \|\vec{\beta}_{A^c}\|^2 + \|\vec{\beta}_A - \hat{\vec{\beta}}_A\|^2$$

a) Case  $p \leq n$

Breiman, Freedman 1983

From the least-squares solution on  $A$  we have (see exercises):

$$\hat{\beta}_A = (X_A^T X_A)^{-1} X_A^T \vec{y}$$

Note that:

$$\text{rank}(X_A^T X_A) \leq \min(n, p) = p$$

Wishart matrix with  $n$  degrees of freedom and covariance matrix  $\Sigma_p$

$\hookrightarrow$  for  $p \leq n$ , it is full rank with high probability.

Let  $\vec{\eta} \equiv \vec{y} - X_A \hat{\beta}_A$  and write:

observe that this is a vector in  $A^c$ :  $\vec{\eta} = X_A \hat{\beta}_A + X_{A^c} \hat{\beta}_{A^c} + \vec{\epsilon} - X_A \hat{\beta}_A$



$$\begin{aligned}
\vec{\beta}_A - \hat{\vec{\beta}}_A &= \vec{\beta}_A - (X_A^T X_A)^{-1} X_A^T (\vec{y} + X_A \vec{\beta}_A) \\
&= \vec{\beta}_A - \underbrace{(X_A^T X_A)^{-1} X_A^T X_A}_{I} \vec{\beta}_A \\
&\quad - (X_A^T X_A)^{-1} X_A^T \vec{y} \\
&= - (X_A^T X_A)^{-1} X_A^T \vec{y}
\end{aligned}$$

Now the square of the  $l_2$ -norm:

$$\begin{aligned}
\|\vec{\beta}_A - \hat{\vec{\beta}}_A\|^2 &= \|(X_A^T X_A)^{-1} X_A^T \vec{y}\|^2 = \\
&= \left( (X_A^T X_A)^{-1} X_A^T \vec{y} \right)^T \left( (X_A^T X_A)^{-1} X_A^T \vec{y} \right) \\
&= \vec{y}^T X_A \left( (X_A^T X_A)^{-1} \right)^T (X_A^T X_A)^{-1} X_A^T \vec{y} = \\
&= \vec{y}^T X_A (X_A^T X_A)^{-2} X_A^T \vec{y}
\end{aligned}$$

Let  $S_A$  be the unique positive definite square root of  $X_A^T X_A$  and define

$$\Psi_A = X_A S_A^{-1} \in \mathbb{R}^{n \times p}$$

Then  $\Psi_A$  is orthonormal:

$$\begin{aligned} \underline{X_A^T X_A} &= (\Psi_A S_A)^T (\Psi_A S_A) = S_A^T \underbrace{\Psi_A^T \Psi_A}_{= \mathbb{I}_p} S_A \\ &= \underline{S_A^2} \end{aligned}$$

### Trace trick

For a matrix  $A \in \mathbb{R}^{p \times p}$  and a vector  $\vec{u} \in \mathbb{R}^d$ , the quantity  $\vec{u}^T A \vec{u}$  is a real number and can be thought as a  $1 \times 1$  matrix. Using the cyclic property of trace:

$$\vec{u}^T A \vec{u} = \text{Tr} \{ \vec{u}^T A \vec{u} \} = \text{Tr} \{ \vec{u} \vec{u}^T A \} = \text{Tr} \{ A \vec{u} \vec{u}^T \}$$

Then:

$$\begin{aligned}\|\vec{\beta}_A - \hat{\vec{\beta}}_A\|^2 &= \vec{\eta}^T \Psi_A S_A S_A^{-4} S_A^T \Psi_A^T \vec{\eta} \\ &= \vec{\eta}^T \Psi_A \underline{S_A^{-2}} \underline{\Psi_A^T \vec{\eta}} = \text{trace trick} \\ &= \vec{v}^T (X_A^T X_A)^{-1} \underline{\vec{v}} \\ &= \text{Tr} \{ \vec{v}^T (X_A^T X_A)^{-1} \vec{v} \} = \\ &= \text{Tr} \{ (X_A^T X_A)^{-1} \vec{v} \vec{v}^T \}\end{aligned}$$

where

$$\vec{v} = \Psi_A^T \vec{\eta} \in \mathbb{R}^p$$

Observe that:

$$\vec{v} = \Psi_A^T (\vec{y} - X \vec{\beta}_A) = \Psi_A^T (X_{A^c} \vec{\beta}_{A^c} + \mu \vec{\epsilon})$$

$$\Rightarrow \mathbb{E}_{X_{A^c}, \vec{\epsilon}} [\vec{v}] = \vec{0}$$

and:

$$\begin{aligned}\vec{v} \vec{v}^T &= \Psi_A^T (X_{A^c} \vec{\beta}_{A^c} + \mu \vec{\epsilon}) (X_{A^c} \vec{\beta}_{A^c} + \mu \vec{\epsilon})^T \Psi_A = \\ &= \Psi_A^T (X_{A^c} \vec{\beta}_{A^c} + \mu \vec{\epsilon}) (\vec{\beta}_{A^c}^T X_{A^c}^T + \mu \vec{\epsilon}^T) \Psi_A = \\ &= \Psi_A^T (X_{A^c} \vec{\beta}_{A^c} \vec{\beta}_{A^c}^T X_{A^c}^T + \mu X_{A^c} \vec{\beta}_{A^c} \vec{\epsilon}^T \\ &\quad + \mu \vec{\epsilon} \vec{\beta}_{A^c}^T X_{A^c}^T + \mu^2 \vec{\epsilon} \vec{\epsilon}^T) \Psi_A\end{aligned}$$

$\Rightarrow$

$$\begin{aligned}\mathbb{E}_{X_{A^c}, \vec{\beta}_{A^c}, \vec{\epsilon}} [\vec{v} \vec{v}^T] &= \\ &= \Psi_A^T \mathbb{E}_{X_{A^c}} [X_{A^c} \vec{\beta}_{A^c} (X_{A^c} \vec{\beta}_{A^c})^T] \Psi_A \\ &\quad + \mathbb{E}_{X_{A^c}} \left[ \Psi_A^T \underbrace{\mathbb{E}_{\vec{\epsilon}} [\vec{\epsilon} \vec{\epsilon}^T]}_{\mathbb{I}_n} \Psi_A \right] \\ &\quad \text{and } \Psi_A^T \Psi_A = \mathbb{I}_p\end{aligned}$$

Observe that:

$$(X \vec{\beta} (X \vec{\beta})^T)_{kb} = (X \vec{\beta})_k (X \vec{\beta})_b'$$

$$= \sum_{l=1}^p x_{kl} \beta_l \sum_{l'=1}^p x_{k'l'} \beta_{l'} \Rightarrow$$

$$\begin{aligned} \mathbb{E} \left[ \left( X \vec{\beta} (X \vec{\beta})^T \right)_{kk'} \right] &= \\ &= \sum_{l=1}^p \beta_l^2 \underbrace{\mathbb{E} [x_{kl} x_{k'l'}]}_{\delta_{kk'}} = \|\vec{\beta}\|^2 \delta_{kk'} \end{aligned}$$

Thus, since

$$\mathbb{E}_{X_{A^c}, \vec{\epsilon}} \left[ X_{A^c} \vec{\beta}_{A^c} (X_{A^c} \vec{\beta}_{A^c})^T \right] = \|\vec{\beta}_{A^c}\|^2 \mathbb{I}_n$$

we have for  $\vec{v}$ :

$$\begin{aligned} \mathbb{E}_{X_{A^c}, \vec{\epsilon}} [\vec{v}] &= \vec{0} \\ \mathbb{E}_{X_{A^c}, \vec{\epsilon}} [\vec{v} \vec{v}^T] &= \left( \|\vec{\beta}_{A^c}\|^2 + \gamma^2 \right) \mathbb{I}_p \end{aligned}$$

Thus, since  $\Psi_A^T \Psi_A = \mathbb{I}_p$  we can

write:

$$\begin{aligned} \mathbb{E}_{X, \vec{y}, \vec{\epsilon}} \left[ \|\vec{\beta}_A - \hat{\vec{\beta}}_A\|^2 \right] &= \mathbb{E}_{X_A, X_{A^c}, \vec{\epsilon}} \left[ \|\vec{\beta}_A - \hat{\vec{\beta}}_A\|^2 \right] \\ &= \text{Tr} \left\{ \mathbb{E}_{X_A} \left[ (X_A^T X_A)^{-1} \right] \mathbb{E}_{X_{A^c}, \vec{\epsilon}} \left[ \vec{v} \vec{v}^T \right] \right\} \\ &= \left( \|\vec{\beta}_{A^c}\|^2 + \eta^2 \right) \text{Tr} \left\{ \mathbb{E}_{X_A} \left[ (X_A^T X_A)^{-1} \right] \mathbb{I}_p \right\} \end{aligned}$$

Observe that  $X_A \in \mathbb{R}^{n \times p}$  is a matrix with  $n$  independent columns sampled from  $\mathcal{N}(\vec{x}_A | \vec{0}; \mathbb{I}_p)$ . Then

$X_A^T X_A$ : Wishart matrix

$$\Downarrow \mathbb{E}_{\vec{x}_A} \left[ X_A^T X_A \right] = n \mathbb{I}_p$$

$$\mathbb{E} \left[ (X_A^T X_A)_{jk} \right] = \sum_{l=1}^n \mathbb{E} \left[ X_A^{lj} X_A^{lk} \right] =$$

$$= \sum_{k=1}^n \delta_{jk} = n \delta_{jk} //$$

Book : Potters and Bouchaud  
 A first course in Random  
 Matrix Theory

The matrix  $(X_A^T X_A)^{-1}$  on the other hand, is an inverse-Wishart matrix, which follows a inverse-Wishart distribution with scale matrix  $\Pi_p$  and  $n$  degrees of freedom.

↳ the distribution is well-defined for  $n > p + 1$

\* When the scale matrix is the identity, diagonal elements of  $(X_A^T X_A)^{-1}$  follow an inverse  $\chi^2$  distribution

: inv- $\chi^2_{n-p+1}$

The first moment of an inverse-Wishart distribution with scale matrix  $\mathbb{I}_p$  and  $n$  degrees of freedom is known to be

$$E_{X_A} \left[ (X_A^T X_A)^{-1} \right] = \frac{\mathbb{I}_p}{n - p - 1}$$

which is well-defined only for  $n > p + 1$

Plugging this result in the expression we had before

$$R_A(\vec{\beta}) = \left( \|\vec{\beta}_{Ac}\|^2 + \eta^2 \right) \left( 1 + \frac{p}{n - p - 1} \right)$$

if  $p \leq n - 2$



# Inverse - Wishart distribution

For a positive definite  $M \in \mathbb{R}^{p \times p}$ , the PDF of the inverse Wishart is

$$f_M(M; S, n) = \frac{(\det S)^{n/2} (\det M)^{-(p+n+1)/2} e^{-\frac{1}{2} \text{Tr}\{SM^{-1}\}}}{2^{pn/2} \Gamma_p\left(\frac{n}{2}\right)}$$

with  $S \in \mathbb{R}^{p \times p}$  positive definite being the scale matrix with, by definition  $S_{jl} \geq 0$ ,  $\forall j, l = 1, \dots, p$ . The number  $n > p + 1$  is called the degrees of freedom. The function  $\Gamma_p(\cdot)$  is the multivariate gamma function.

The first moment is

$$E_M[M] = \frac{S}{n - p - 1}; \quad n > p + 1$$

b) Case  $p \geq n$

Remembering that  $\hat{\beta}_A = X_A^T y$  and writing the Moore-Penrose inverse as

$$\underline{X_A^{\dagger} = X_A^T (X_A X_A^T)^{\dagger}}$$

Property of the Moore-Penrose inverse when  $X_A X_A^T$  and  $X_A^T X_A$  are symmetric

and defining:

$$\vec{v} \equiv y - X_A \hat{\beta}_A$$

we write:

$$\begin{aligned} \hat{\beta}_A - \underline{\hat{\beta}}_A &= \hat{\beta}_A - X_A^{\dagger} y = \\ &= \hat{\beta}_A - X_A^T (X_A X_A^T)^{\dagger} (\vec{v} + X_A \hat{\beta}_A) \\ &= \hat{\beta}_A - X_A^T (X_A X_A^T)^{\dagger} X_A \hat{\beta}_A \end{aligned}$$

$$\begin{aligned}
& - X_A^T (X_A X_A^T)^{\dagger} \vec{y} \\
& = \left( \mathbb{I}_p - X_A^T (X_A X_A^T)^{\dagger} X_A \right) \vec{\beta}_A \\
& \quad - X_A^T (X_A X_A^T)^{\dagger} \vec{y} \\
& = \underline{\left( \mathbb{I}_p - \mathcal{P}_{X_A} \right) \vec{\beta}_A} - \underline{X_A^T (X_A X_A^T)^{\dagger} \vec{y}}
\end{aligned}$$

where we have defined the projection matrix onto the row space of  $X_A$ :

$$\mathcal{P}_{X_A} \equiv X_A^T (X_A X_A^T)^{\dagger} X_A$$

$\Rightarrow \left( \mathbb{I}_p - \mathcal{P}_{X_A} \right) \vec{\beta}_A$  : orthogonal projection of  $\vec{\beta}_A$  onto the (kernel) null space of  $X_A$

$$\underline{X_A^T (X_A X_A^T)^{\dagger} \vec{y}}$$

$\Rightarrow X_A^T (X_A X_A^T)^{\dagger} \vec{y}$  : vector on the  
new space of  
 $X_A$

These two vectors are then  
orthogonal:

$$\|\vec{\beta}_A - \hat{\vec{\beta}}_A\|^2 =$$

$$= \underbrace{\|(\mathbb{I}_p - \mathcal{P}_{\hat{X}_A}) \vec{\beta}_A\|^2} + \underbrace{\|X_A^T (X_A X_A^T)^{\dagger} \vec{y}\|^2}$$

First term

$$\|(\mathbb{I}_p - \mathcal{P}_{\hat{X}_A}) \vec{\beta}_A\|^2$$

$$= \vec{\beta}_A^T (\mathbb{I}_p - \mathcal{P}_{\hat{X}_A}^T) (\mathbb{I}_p - \mathcal{P}_{\hat{X}_A}) \vec{\beta}_A =$$

$$= \vec{\beta}_A^T \vec{\beta}_A - \vec{\beta}_A^T \mathcal{P}_{X_A} \vec{\beta}_A - \vec{\beta}_A^T \mathcal{P}_{X_A}^T \vec{\beta}_A + \vec{\beta}_A^T \mathcal{P}_{X_A}^T \mathcal{P}_{X_A} \vec{\beta}_A$$

Observe that:

$$\begin{aligned} P_{X_A}^T &= \left( X_A^T (X_A X_A^T)^{\dagger} X_A \right)^T = X_A^T \left( X_A^T (X_A X_A^T)^{\dagger} \right)^T \\ &= X_A^T \left( (X_A X_A^T)^{\dagger} \right)^T X_A = X_A^T (X_A X_A^T)^{\dagger} X_A \\ &= P_{X_A} \end{aligned}$$

Also note that:

$$P_{X_A}^2 = X_A^T \underbrace{(X_A X_A^T)^{\dagger}}_{n \times n \text{ full rank matrix with high probability}} X_A X_A^T (X_A X_A^T)^{\dagger} X_A$$

$n \times n$  full rank matrix with high probability, as  $\vec{X}_A$  is full row rank with high probability: has "almost"  $n$  linearly independent rows (with high probability)

$$\Rightarrow (X_A X_A^T)^{\dagger} X_A X_A^T \approx I_n$$

$$= X_A^T (X_A X_A^T)^{\dagger} X_A = P_{X_A}$$

Then we have:

$$\begin{aligned} & \|(\mathbb{I}_p - \mathcal{P}_{X_A}) \vec{\beta}_A\|^2 = \\ &= \vec{\beta}_A^T \vec{\beta}_A - 2 \vec{\beta}_A^T \mathcal{P}_{X_A}^T \mathcal{P}_{X_A} \vec{\beta}_A + \vec{\beta}_A^T \mathcal{P}_{X_A}^T \mathcal{P}_{X_A} \vec{\beta}_A \\ &= \vec{\beta}_A^T \vec{\beta}_A - (\mathcal{P}_{X_A} \vec{\beta}_A)^T (\mathcal{P}_{X_A} \vec{\beta}_A) \end{aligned}$$

$\Rightarrow$

$$\begin{aligned} & \|(\mathbb{I}_p - \mathcal{P}_{X_A}) \vec{\beta}_A\|^2 = \\ &= \|\vec{\beta}_A\|^2 - \|\mathcal{P}_{X_A} \vec{\beta}_A\|^2 \end{aligned}$$

By rotation symmetry of the standard normal distribution we have:

$$E_{X_A, \vec{\beta}_A} [\|\mathcal{P}_{X_A} \vec{\beta}_A\|^2] = \frac{r}{p} \|\vec{\beta}_A\|^2$$

\* at the end  
there is a remark on this

⇒

$$\begin{aligned} E_{X_A} \left[ \left\| (\mathbb{I}_p - \mathcal{P}_{X_A}) \vec{\beta}_A \right\|^2 \right] \\ = \left\| \vec{\beta}_A \right\|^2 \left( 1 - \frac{n}{p} \right) \end{aligned}$$

Second term

$$\begin{aligned} \left\| X_A^T (X_A X_A^T)^{\dagger} \vec{\eta} \right\|^2 &= \\ &= \text{Tr} \left\{ (X_A^T (X_A X_A^T)^{\dagger} \vec{\eta})^T X_A^T (X_A X_A^T)^{\dagger} \vec{\eta} \right\} = \\ &= \text{Tr} \left\{ \vec{\eta}^T (X_A^T (X_A X_A^T)^{\dagger})^T X_A^T (X_A X_A^T)^{\dagger} \vec{\eta} \right\} = \\ &= \text{Tr} \left\{ ((X_A X_A^T)^{\dagger})^T X_A X_A^T (X_A X_A^T)^{\dagger} \vec{\eta} \vec{\eta}^T \right\} = \\ &= \text{Tr} \left\{ (X_A X_A^T)^{\dagger} X_A X_A^T (X_A X_A^T)^{\dagger} \vec{\eta} \vec{\eta}^T \right\} \\ &\quad \underbrace{\hspace{10em}}_{\text{almost surely } \mathbb{I}_n \text{ because}} \end{aligned}$$

$X_A X_A^T \in \mathbb{R}^{n \times n}$  ( $n < p$ ) is almost  
surely invertible  
(with high probability)

$$= \text{Tr} \{ (X_A X_A^T)^+ \vec{\eta} \vec{\eta}^T \}$$

Remember that:

$$\vec{\eta} = \vec{y} - \underline{X_A \vec{\beta}_A}$$

$$\underline{X_A \vec{\beta}_A} + X_{A^c} \vec{\beta}_{A^c} + y \vec{e}$$

is a vector in  $A^c$ :

$$\vec{\eta} = X_{A^c} \vec{\beta}_{A^c} + y \vec{e} \in \mathbb{R}^n$$

Then since  $X_A \vec{\beta}_A$  and  $X_{A^c} \vec{\beta}_{A^c} + y \vec{e}$   
are uncorrelated we have



$$\mathbb{E}_{x, \vec{y}} \left\{ \left\| X_A^T (X_A X_A^T)^{\dagger} \vec{y} \right\|^2 \right\} =$$

$$= \text{Tr} \left\{ \mathbb{E}_{x, y} \left[ (X_A X_A^T)^{\dagger} \right] \mathbb{E}_{x, y} \left[ \vec{y} \vec{y}^T \right] \right\} =$$

$$= \text{Tr} \left\{ \underbrace{\mathbb{E}_{X_A} \left[ (X_A X_A^T)^{\dagger} \right]}_{\text{blue underline}} \underbrace{\mathbb{E}_{X_{A^c}} \mathbb{E}_{\vec{e}} \left[ \vec{y} \vec{y}^T \right]}_{\text{blue underline}} \right\}$$

Note that

$$\begin{aligned} \vec{y} \vec{y}^T &= (X_{A^c} \vec{\beta}_{A^c} + \mu \vec{e}) (\vec{\beta}_{A^c}^T X_{A^c}^T + \mu \vec{e}^T) \\ &= \|\vec{\beta}_{A^c}\|^2 X_{A^c} X_{A^c}^T + \mu X_{A^c} \vec{\beta}_{A^c} \vec{e}^T \\ &\quad + \mu \vec{e} \vec{\beta}_{A^c}^T X_{A^c}^T + \mu^2 \vec{e} \vec{e}^T \end{aligned}$$

$$\Rightarrow \underbrace{\mathbb{E}_{x, y} \mathbb{E}_{\vec{e}} \left[ \vec{y} \vec{y}^T \right]}_{\text{blue underline}} = (\|\vec{\beta}_{A^c}\|^2 + \mu^2) \mathbb{I}_n$$

Let us now define :

$$\underline{\underline{\mathbb{I}_A \equiv (X_A X_A^T)^T}}$$

Observe that  $X_A \in \mathbb{R}^{n \times p}$  is a matrix with  $n$  independent columns sampled from  $\mathcal{N}^p(\vec{x}_A | \vec{0}; \mathbb{1}_p)$ . Then

$X_A X_A^T$  : Wishart matrix

$$\Downarrow \underline{\underline{\mathbb{E}_{X_A} [X_A X_A^T] = p \mathbb{I}_n}}$$

$$\begin{aligned} \mathbb{E} [(X_A X_A^T)_{jk}] &= \sum_{l=1}^p \mathbb{E} [X_A^{jl} X_A^{kl}] = \\ &= \sum_{l=1}^p \delta_{jk} = p \delta_{jk} // \end{aligned}$$

Book : Potters and Bouchaud  
A first course in Random  
Matrix Theory

# Inverse - Wishart distribution

For a positive definite  $M \in \mathbb{R}^{n \times n}$ , the PDF of the inverse Wishart is

$$f_M(M; S, p) = \frac{(\det S)^{p/2}}{2^{pn/2} \Gamma_n\left(\frac{p}{2}\right)} (\det M)^{-(p+n+1)/2} e^{-\frac{1}{2} \text{Tr}\{SM^{-1}\}}$$

with  $S \in \mathbb{R}^{n \times n}$  positive definite being the scale matrix with, by definition  $S_{jl} \geq 0$ ,  $\forall j, l = 1, \dots, n$ . The number  $p > n+1$  is called the degrees of freedom. The function  $\Gamma_n(\cdot)$  is the multivariate gamma function.

The first moment is

$$\mathbb{E}_M[M] = \frac{S}{p - n - 1}; \quad p > n + 1$$

The matrix  $\mathbb{I}_A^{-1}$ , on the other hand, is an inverse-Wishart matrix, which follows a inverse-Wishart distribution with scale matrix  $\mathbb{I}_n$  and  $p$  degrees of freedom.

↳ the distribution is well-defined for  $p > n+1$

\* When the scale matrix is the identity, diagonal elements of  $\mathbb{I}_A$  follow an inverse  $\chi^2$  distribution

$$: \text{inv-}\chi^2_{p-n+1}$$

The first moment of an inverse-Wishart distribution with scale matrix  $\mathbb{I}_n$  and  $p$  degrees of freedom is known to be

$$\mathbb{E}_{X_A} [\mathbb{I}_A] = \frac{\mathbb{I}_n}{p - n - 1}$$

which is well-defined only for  $p > n+1$ .

For  $p = n+1$ , an extrapolation of the result above makes the first moment infinite, which can be interpreted as the PDF going to zero, as one can write  $e^{-Tn^2-3}$  in terms of the first moment.

The case  $p = n$  can also be interpreted as the expectation going to infinite in order to send the PDF to zero.

↳ It is also consistent with Breiman, Freedman (1983)

Therefore for  $p > n+1$ :

$$\mathbb{E}_{X, \vec{y}} \left\{ \left\| X_A^T (X_A X_A^T)^{-1} \vec{y} \right\|^2 \right\} =$$

$$\text{Tr} \left\{ \frac{\mathbb{I}_n}{p-n-1} \left( \|\vec{\beta}_{A^c}\|^2 + \eta^2 \right) \mathbb{I}_n \right\}$$

$$= \frac{n}{p-n-1} (\|\vec{\beta}_{A^c}\|^2 + \eta^2)$$

Summing the **first term** with the **second term** we finally obtain:

$$\mathcal{R}_A(\vec{\beta}) =$$

$$\begin{cases} +\infty ; & \text{if } n-1 \leq p \leq n+1 \\ \|\vec{\beta}_A\|^2 \left(1 - \frac{n}{p}\right) + (\|\vec{\beta}_{A^c}\|^2 + \eta^2) \left(1 + \frac{n}{p-n-1}\right) & \text{if } p \geq n+2 \end{cases}$$

Remark on:

$$\mathbb{E}_{X, \vec{y}} \left[ \left\| \mathcal{P}_{X_A} \vec{\beta}_A \right\|^2 \right] = \frac{n}{p} \left\| \vec{\beta}_A \right\|^2$$

Observe that

$$\begin{aligned} \left\| \mathcal{P}_{X_A} \vec{\beta}_A \right\|^2 &= \left( \mathcal{P}_{X_A} \vec{\beta}_A \right)^T \left( \mathcal{P}_{X_A} \vec{\beta}_A \right) = \\ &= \vec{\beta}_A^T \mathcal{P}_{X_A}^T \mathcal{P}_{X_A} \vec{\beta}_A = \vec{\beta}_A^T \mathcal{P}_{X_A}^2 \vec{\beta}_A = \\ &= \vec{\beta}_A^T \mathcal{P}_{X_A} \vec{\beta}_A = \text{Tr} \left\{ \vec{\beta}_A^T \mathcal{P}_{X_A} \vec{\beta}_A \right\} = \\ &= \text{Tr} \left\{ \vec{\beta}_A \vec{\beta}_A^T \mathcal{P}_{X_A} \right\} \end{aligned}$$

$$\begin{aligned} \Rightarrow \mathbb{E}_{X_A} \left[ \left\| \mathcal{P}_{X_A} \vec{\beta}_A \right\|^2 \right] &= \\ &= \text{Tr} \left\{ \vec{\beta}_A \vec{\beta}_A^T \mathbb{E}_{X_A} \left[ \mathcal{P}_{X_A} \right] \right\} \end{aligned}$$

The matrix  $\mathcal{P}_{X_A}$  can be viewed as a projection matrix, that projects vectors

in  $\mathbb{R}^p$  onto the column space of  $X_A$ .

Given that the rows of  $X_A$  are iid  $\mathcal{N}(\vec{x}_i | \vec{0}_p, \mathbb{I}_p)$ , none of the  $p$  directions should be preferred. Then it's reasonable to expect:

$$\mathbb{E}_{X_A} [P_{X_A}] \propto \mathbb{I}_p$$

isotropic  
rotational symmetry

Since  $P_{X_A}$  is a projection matrix of rank  $n$  with high probability, its trace (the sum of its eigenvalues, which are either 0 or 1; remember that  $P_{X_A} = P_{X_A}^2$ ) should be  $n$ , with high probability. Then, it's reasonable to expect:

$$\text{Tr} \left\{ \mathbb{E}_{X_A} [P_{X_A}] \right\} = n$$

Thus assuming that  $\mathbb{E}_{X_A} [P_{X_A}]$  must be isotropic and proportional to



$\frac{\text{Tr} \mathbb{P}}{p}$  the scale factor that ensures  $\text{Tr} \{ E_{X_A} [P_{X_A}] \} = n$  is  $n/p$ . Then:

$$E_{X_A} [P_{X_A}] = \frac{n}{p} \mathbb{P}$$

The projection  $P_{X_A}$  essentially distributes the effect of the  $n$  dimensions uniformly (on average) across the  $p$  components.

Finally:

$$\begin{aligned} E_{X_A} [ \| P_{X_A} \vec{\beta}_A \|^2 ] &= \text{Tr} \{ \vec{\beta}_A \vec{\beta}_A^T E_{X_A} [P_{X_A}] \} \\ &= \text{Tr} \left\{ \vec{\beta}_A \vec{\beta}_A^T \mathbb{P} \frac{n}{p} \right\} \\ &= \frac{n}{p} \text{Tr} \{ \vec{\beta}_A \vec{\beta}_A^T \} \\ &= \frac{n}{p} \sum_{j=1}^p \beta_{A,j}^2 = \frac{n}{p} \| \vec{\beta}_A \|^2 \end{aligned}$$