

Chapter 11: Graphical Models for High-Dimensional Data:

11.1 $S^{d \times d}$: Set of symmetric matrices

$S_+^{d \times d}$: Cone of symmetric and strictly positive matrices

$$F: S^{d \times d} \rightarrow \mathbb{R} \text{ given by } F(\Theta) = \begin{cases} -\sum_{j=1}^d \log \gamma_j(\Theta) & \text{if } \Theta \in S_+^{d \times d} \\ +\infty & \text{otherwise} \end{cases}$$

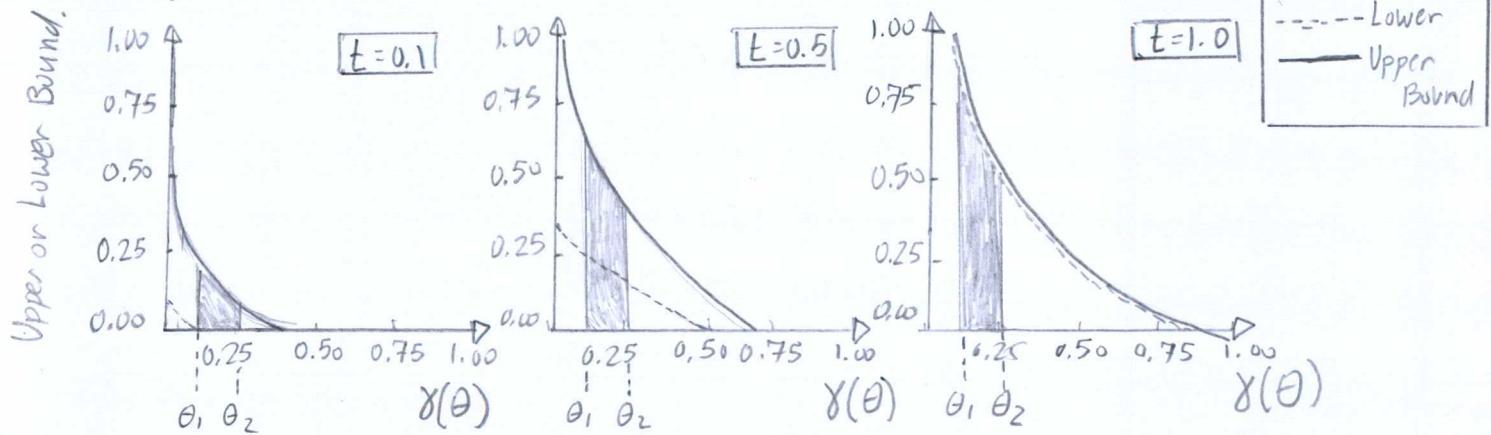
where $\gamma_j(\Theta) > 0$ are eigenvalues

of Θ

a) F is strictly convex in the $S_+^{d \times d}$

domain by Jensen's Inequality: $F(tX_1 + (1-t)X_2) \leq tF(X_1) + (1-t)F(X_2)$

$$-\sum_{j=1}^d \log \gamma(t\theta_1 + (1-t)\theta_2) \leq t \left(-\sum_{j=1}^d \log \gamma(\theta_1) \right) + (1-t) \left(-\sum_{j=1}^d \log \gamma(\theta_2) \right)$$



$$\begin{aligned} b) \nabla F(\Theta) &= \nabla \left(-\sum_{j=1}^d \log \gamma_j(\Theta) \right) = -\frac{1}{\Theta} = \Theta^{-1} \\ &= -1/\Theta \\ &= -\Theta^{-1} \end{aligned}$$

$$c) \nabla F^2(\Theta) = \nabla \left(-\sum_{j=1}^d \log \gamma_j(\Theta) \right)^2$$

$$= 1/\Theta^2$$

$$= \Theta^{-1} \otimes \Theta^{-1}$$

$$11.2 \quad \hat{\Theta}_{MLE} = \begin{cases} \Sigma^{-1} & \text{if } \hat{\Sigma} > 0 \\ \text{not defined} & \text{otherwise} \end{cases} \quad \text{where } \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X_i X_i^T$$

$$\text{Gaussian: } f(X|0, \Sigma) = \frac{1}{\sqrt{2\pi\Sigma}} e^{-\frac{X^T X}{2\Sigma}}$$

$$\begin{aligned} \frac{df(X|0, \Sigma)}{d\Sigma} &= \frac{d}{d\Sigma} \left(\frac{1}{\sqrt{2\pi\Sigma}} e^{-\frac{X^T X}{2\Sigma}} \right) \\ &= -\frac{1}{2} \left(\frac{1}{\sqrt{2\pi\Sigma\Sigma}} \right) e^{-\frac{X^T X}{2\Sigma}} + \frac{1}{2} \left(\frac{1}{\sqrt{2\pi\Sigma\Sigma^2}} \right) e^{-\frac{X^T X}{2\Sigma}} \end{aligned}$$

$$= 0$$

$$\Sigma^* = \boxed{1\Theta_{xx}}$$

$$\hat{\Theta}_{MLE} = \Sigma^{-1}$$

$$11.3. \quad X \in \mathbb{R}^d, \quad \Sigma^*, \quad Z^* = (X_3 | X_{1333})$$

$$a) \mathbb{E}[Z] = \mathbb{E}[X_3 | X_{1333}]$$

$$= W_3 + \theta \cdot X_{1333} - X_3$$

$$= 0$$

$$X_3 = \langle \theta, X_{1333} \rangle + W_3$$

$$b) \mathbb{E}[(X_{1333} - \theta X_3) X_{1333}^T] = \mathbb{E}[X_{1333} \circ X_{1333}] - \theta \mathbb{E}[X_{1333} \circ X_3]$$

$$= \Sigma_{1333, 1333} - \theta \Sigma_{1333, 3}$$

$$= 0$$

$$\theta = \frac{\Sigma_{1333, 3}}{\Sigma_{1333, 1333}}$$

c) The "just" in the problem relates a covariance matrix to zero-mean Gaussians with zero-off-diagonal elements e.g. no variance between the x-, and y-directions.

The matrix is a "diagonally dominant" covariance.

$$a_{ii} > \sum_{j \neq i} |a_{ij}| = a_{ii} + a_{iz} + \dots + a_{ij}$$

$$\begin{aligned}\Sigma &= \frac{\Sigma_{21} \Sigma_{12}}{\Sigma_{11} \circ \Sigma_{11}} \\ &= \begin{bmatrix} I & -A_{11}^{-1} A_{12} \\ -A_{21} & A_{22} \end{bmatrix} \\ &\stackrel{1/2}{=} \left(\frac{\Sigma_{21} \Sigma_{12}}{\Sigma_{11} \circ \Sigma_{11}} \right)^{1/2} \\ &= \begin{bmatrix} I & 0 \\ -\frac{A_{21}}{A_{11}} & I \end{bmatrix} \begin{bmatrix} I & -\frac{A_{12}}{A_{11}} \\ 0 & I \end{bmatrix}\end{aligned}$$

$$\begin{aligned}\hat{\Theta}_{jk}^* &= \Sigma^{1/2} \cdot A \circ \Sigma^{1/2} \\ &= \begin{bmatrix} I & 0 \\ -\frac{A_{21}}{A_{11}} & I \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} I & -\frac{A_{12}}{A_{11}} \\ 0 & I \end{bmatrix} \\ &= \begin{bmatrix} A_{11} & 0 \\ 0 & A_{22} - A_{21} \circ A_{11}^{-1} A_{12} \end{bmatrix} \\ &= \underbrace{\begin{bmatrix} A_{11} & 0 \\ 0 & A_{22} - A_{21} \circ A_{11}^{-1} A_{12} \end{bmatrix}}_{\text{"Strictly Diagonal"}},\end{aligned}$$

$$\text{II.4 } \hat{\Theta} = \arg \min_{\Theta \in \mathbb{R}^{d \times d}} \{ \|\Theta\|_1 \} \text{ such that } \|\hat{\Sigma} \Theta - I_d\|_{\max} \leq \lambda_n$$

$$\hat{T} = \arg \min_{T \in \mathbb{R}^{d \times d}} \{ \|T\|_1 \} \text{ such that } \|\hat{\Sigma} T - e_j\|_{\max} \leq \lambda_n$$

where $e_j \in \mathbb{R}^d$ is the j^{th} canonical vector

$$\text{a) } \hat{\Theta} = \arg \min_{\Theta} \left\{ \frac{1}{2} \langle \hat{\Theta}^2, \Sigma \rangle - \text{tr}(\Theta) + \lambda \|\Theta\|_1 \right\} = 0$$

$$\hat{\lambda}^* = \frac{1}{2} (\Theta \hat{\Sigma} + \hat{\Sigma} \Theta) - I$$

$$\hat{T} = \arg \min_{T} \left\{ \frac{1}{2} \langle \hat{T}^2, \Sigma \rangle - e_j \text{tr}(T) + \lambda \|\Theta\|_1 \right\} = 0$$

$$\hat{\lambda}^* = \frac{1}{2} (T \Sigma + \Sigma T) - e_j$$

An optimal solution depends on the j^{th} vector

b) When $\lambda \geq \|\theta^*\|_1, \|\hat{\Sigma} - \Sigma\|_{\max}$

$$\begin{aligned}\|\hat{T}_j\|_1 &= \operatorname{argmax}_{T_j} \left\{ \frac{1}{2} T \Sigma T^T - e T_j + \lambda_n \|T_j\|_1 \right\} \\ &= \operatorname{argmin}_{T_j} \left\{ \frac{1}{2} T \Sigma T^T - e T_j + \|\theta^*\|_1, \|\hat{\Sigma} - \Sigma\|_{\max} \|T_j\|_1 \right\} \\ &\leq \operatorname{argmin}_{\theta} \left\{ \frac{1}{2} \theta \Sigma \theta^T - I + \|\theta^*\|_1, \|\hat{\Sigma} - \Sigma\|_{\max} \|\theta\|_1 \right\} \\ &\leq \|\theta_j\|_1\end{aligned}$$

c) $\max_{j=1, \dots, d} \Sigma_{jj} \leq 1 \Rightarrow \|\hat{\Sigma} - \Sigma^*\|_1 / \|\theta^*\|_1 \leq \lambda$

$$\|\hat{\Sigma} - \Sigma^*\|_1 \leq \frac{\lambda_{\max}}{\|\theta^*\|_1}$$

where λ_{\max} is the maximum eigenvalue

Two or three papers support the bound by

Corollary or Propositions:

- 1) Eigenvalue bounds are necessary
- 2) A data dimension relationship: $\sigma^2 \log d = O(n)$
- 3) Exponential tails on the graph

Outcomes from the model in part d.

- A) Continuous, increasing and a concave plot from 0 to inf.
- B) Zero at the origin, singularity.
- C) Constants in the coefficient for fits at all inputs.

d) If $\lambda = C_1 \sqrt{\frac{\log d}{n}}$

$$\|\hat{\theta} - \theta^*\|_1 \leq C_1 \|\theta^*\|_1 \sqrt{\frac{\log d}{n}}$$

$$\begin{aligned}
 & \text{II.5. (Equation II.37)} \\
 \hat{\theta}_{j+}^* &= \underset{\theta_{j+}}{\operatorname{argmin}} \left\{ -\frac{1}{n} \sum_{i=1}^n \log p_{\theta_j^*} \{X_i | X_{i,j+}\} + \lambda \sum_{k \in V \setminus \{j+1\}} \|\theta_{jk}^*\|_1 \right\} \\
 &= \underset{\theta_{j+}}{\operatorname{argmin}} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n v_{ij} g_i \right\|_2^2 + \lambda \sum \|\theta_{jk}\|_1 \right\} \\
 &= \underset{\theta_{j+}}{\operatorname{argmin}} \underbrace{\left\{ \left\| \frac{1}{n} \sum_{i=1}^n \{X_i - X_{i,j+} \theta\} \right\|^2 \right\}}_{L_n(\theta_{j+}; X_i, X_{i,j+})} + \lambda n \|\theta\|_1
 \end{aligned}$$

The regression term is a Lagrangian from
Equation 9.66

$$\begin{aligned}
 & \text{II.6. (Equation II.32)} p_{\theta}(X_1, \dots, X_d) \propto \left\{ \sum_{j \in V} \Phi_j(X_j; \theta_j^*) + \sum_{j,k} \Phi_{jk}(X_j, X_k; \theta_{jk}^*) \right\} \\
 & \text{where } \Phi_j(X_j; \theta_j^*) = \theta_j^* X_j \\
 & \Phi_{jk}(X_j, X_k; \theta_{jk}^*) = \theta_{jk}^* X_j X_k
 \end{aligned}$$

$$\theta_{jk}^* = \begin{cases} \theta_{jk}^* X_j X_k & (j, k) \in E \text{ "points in a graph's vector set"} \\ 0 & (j, k) \notin E \text{ "points not in a graph's vector set"} \end{cases}$$

$$\text{II.7 } X \in \{-1, 1\}^d ; p_{\theta}(x_1, \dots, x_d) = \exp \left\{ \sum_{(j, k) \in E} \theta_{jk} \cdot x_j \cdot x_k - \phi(\theta) \right\}$$

$$\begin{aligned}
 \text{a) For each edge } (j, k) \in E, \frac{d\phi(\theta)}{d\theta_{jk}} &= \frac{1}{d} \sum_{(i, k) \in E} x_i x_i \\
 &= \mathbb{E}[x_j x_k]
 \end{aligned}$$

$$\text{b) } P(X_1=1, \theta) = \exp \left\{ \sum_{(j, k) \in E} \theta_{jk} \right\}$$

$$P(X_1=-1, \theta) = \exp \left\{ -\sum_{(j, k) \in E} \theta_{jk} \right\}$$

$$\begin{aligned}
 P(X_k | X_{-k}) &\equiv \frac{P(X_k = 1)}{P(X_k = 1) + P(X_k = -1)} \\
 &= \frac{\exp \left\{ \sum_{(j,k) \in E} \theta_j X_j \right\}}{\exp \left\{ \sum_{(j,k) \in E} \theta_j X_j \right\} + \exp \left\{ -\sum_{(j,k) \in E} \theta_j X_j \right\}} \\
 &= \frac{\exp \left\{ \sum_{(j,k) \in E} \theta_j X_j \right\}}{\exp \left\{ \sum_{(j,k) \in E} \theta_j X_j \right\} + 1} \left(\frac{1}{\exp \left\{ -\sum_{(j,k) \in E} \theta_j X_j \right\}} \right) \\
 &= \frac{\exp \left(2 \sum_{(j,k) \in E} \theta_j X_j \right)}{\exp \left(\sum_{(j,k) \in E} \theta_j X_j \right) + 1}
 \end{aligned}$$

$$\frac{P(X_k | X_{-k})}{1 - P(X_k | X_{-k})} = \exp \left(\sum_{(j,k) \in E} \theta_j X_j \right)$$

$$\sum_{(j,k) \in E} \theta_j X_j = \log \frac{P(X_k | X_{-k})}{1 - P(X_k | X_{-k})}$$

$$\text{I.I.g. } X = (X_1, \dots, X_d) = N(X | 0, \sigma^2)$$

$$V = N(V | 0, \sigma^2 I_d)$$

$$Z = X + V$$

$$\begin{aligned}
 \text{If } \sigma^2 \|\theta^*\|_2 \leq 1, \text{ then } Z &= X + V \\
 &= N(X | 0, \sigma_x^2 I_d) + N(V | 0, \sigma_y^2 I_d)
 \end{aligned}$$

$$= N(X | 0, (\sigma_x^2 + \sigma_y^2) I_d)$$

$$\frac{1}{\sigma_z^2} = \frac{1}{\sigma_x^2} + \frac{1}{\sigma_y^2}$$

"Expand out in powers"

$$1) e^{-x^2/2\sigma_z^2} = \sum_{n=0}^{\infty} \frac{1}{n!} \left(-\frac{x^2}{2\sigma_z^2} \right)^n$$

$$= 1 - \frac{x^2}{2\sigma_z^2} + \frac{x^4}{8\sigma_z^4} - \frac{x^6}{48\sigma_z^6} + \dots$$

$$= 1 - \left(\frac{1}{2}\right) \left(\frac{\sigma_x^2 \sigma_y^2}{\sigma_x^2 + \sigma_y^2}\right) x^2 + \left(\frac{1}{8}\right) \left(\frac{\sigma_x^4 \sigma_y^4}{(\sigma_x^2 + \sigma_y^2)^2}\right) x^4 - \dots$$

2) At $x = 1$, $\sigma_z^2 = \begin{pmatrix} 1 & \sqrt{\frac{\sigma_x^2 \sigma_y^2}{2(\sigma_x^2 + \sigma_y^2)}} \\ \sqrt{\frac{\sigma_x^2 \sigma_y^2}{2(\sigma_x^2 + \sigma_y^2)}} & 1 \end{pmatrix}$

Cholesky
Decomposition.

$$= \begin{pmatrix} \sqrt{1} & 0 \\ \sqrt{\frac{\sigma_x^2 \sigma_y^2}{2(\sigma_x^2 + \sigma_y^2)}} & \sqrt{1 - \frac{\sigma_x^2 \sigma_y^2}{2(\sigma_x^2 + \sigma_y^2)}} \end{pmatrix} \begin{pmatrix} \sqrt{1} & \sqrt{\frac{\sigma_x^2 \sigma_y^2}{2(\sigma_x^2 + \sigma_y^2)}} \\ 0 & \sqrt{1 - \frac{\sigma_x^2 \sigma_y^2}{2(\sigma_x^2 + \sigma_y^2)}} \end{pmatrix}$$
$$= \sigma_z^0 \sigma_z^1$$

In any case, whatever the expansion, the expansion represents an accumulative deviation from each graph vertex. The standard deviation is amplitude in a Gaussian function and spread too.

11.9. (Equation 11.44) $\hat{\theta} = \underset{\theta \in S_T^{d \times d}}{\operatorname{argmin}} \{ \langle \langle \theta, \hat{T} \rangle \rangle - \log \det \theta + \lambda \| \theta \|_1 \}$

a) $\Sigma_{\text{cov}}(X), \lambda_n > \|\hat{T} - \Sigma_X\|_{\max}$

$$\hat{\theta} = \underset{\theta \in S_T^{d \times d}}{\operatorname{argmin}} \{ \langle \langle \theta, T \rangle \rangle - \log \det \theta + \|\hat{T} - \Sigma_X\|_{\max} \| \theta \|_1 \}$$
$$= 0$$

$$\theta^* = \frac{1}{T + \|T - \Sigma\|}$$

$$\tilde{\theta} = \log(T + \|T - \Sigma\|) + 1$$

b) When λ has no bound, $\tilde{\theta}$ has a singularity at zero e.g. no value.

$$11.10. y = X\theta^* + w \quad ; \text{rowspan}(X) = N(x|0, \Sigma)$$

$$z = X + V \quad ; \text{rowspan}(V) = N(x|0, \sigma^2 I_d)$$

$$\tilde{\theta} = \arg \min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{2n} \|y - z\theta\|_2^2 \right\}$$

Two inconsistencies exist in the problem:

i) Z 's variance is potentially imaginary.

$$Z = X + V$$

$$= N(x|0, \Sigma) + N(x|0, \sigma^2 I_d)$$

$$= N(x|0, \Sigma + \sigma^2 I_d)$$

$$\text{From 11.8, } \frac{1}{\sigma_Z^2} = \frac{1}{\Sigma} + \frac{1}{\sigma^2 I_d}$$

Taylor Expansion to

$$\sigma_Z^2 = \frac{\Sigma \sigma^2 I_d}{\Sigma + \sigma^2 I_d}$$

Second-order

$$= \begin{pmatrix} 1 & \sqrt{\frac{\Sigma \sigma^2 I_d}{\Sigma + \sigma^2 I_d}} \\ \sqrt{\frac{\Sigma \sigma^2 I_d}{\Sigma + \sigma^2 I_d}} & 1 \end{pmatrix}$$

Cholesky

Decomposition

$$= \begin{pmatrix} 1 & 0 & \checkmark & 1 & 0 \\ \sqrt{\frac{\Sigma \sigma^2 I_d}{\Sigma + \sigma^2 I_d}} & \sqrt{1 - \frac{\Sigma \sigma^2 I_d}{\Sigma + \sigma^2 I_d}} & & 0 & \sqrt{1 - \frac{\Sigma \sigma^2 I_d}{\Sigma + \sigma^2 I_d}} \end{pmatrix}$$

The root $\sqrt{1 - \frac{\sum \sigma^2 \text{Id}}{\sum t \sigma^2 \text{Id}}}$ needs a bound
in the problem, such as $\frac{\sum \sigma^2 \text{Id}}{\sum t \sigma^2 \text{Id}} \leq 1$.

2) A rare case when $\text{rowspan}(V)$ draws from $\text{rowspan}(X)$ as in this problem - a zero determinant in the denominator of $\arg\min\{\cdot\}$.

$$\hat{\theta} = \arg\min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 \right\}$$

$$\theta^* = \arg\min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{2n} \|X\theta + w - (X + V)\theta\|_2^2 \right\}$$

$$= \frac{1}{n} \|w - V\theta\|_2 (\text{---})$$

$$= 0$$

$$\theta^* = \frac{w}{V} \quad \text{The denominator has a possible zero division}$$

A solution involves a "regularizer" term, λ .

$$\hat{\theta} = \arg\min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 + \lambda \|\theta\|_1 \right\}$$

11.11 What is a high-dimensional dataset?
n-rows (measurements) and d-columns (data points)

Where $n < d$,

Although, the upper limit $(\sqrt{\frac{n}{\log d}})$ bounds the maximum size. See the following relationship

Format:

$$\min_{\|\theta\|_1 \leq \sqrt{\frac{n}{\log d}}} \left\{ f_n(\theta) + \lambda \|\theta\|_1 \right\}$$

Functional limit

- avoids trivial solutions

Regularization term

- prevents zero division

Three examples explain the functional limit below.

Notes: Papers online study beyond $\sqrt{\frac{n}{\log d}}$ in

datasets. The documents suggest autoregressive patterns in outputs, such as overfit tails, oscillatory data, and "waves" in the image's fit.

Example #1: Additive Noise

A function category, $\min_{\|\theta\|_1 \leq \sqrt{\frac{n}{\log d}}} \left\{ \frac{1}{2} \theta^T \hat{T} \theta - \langle \hat{y}, \theta \rangle + \lambda \|\theta\|_1 \right\}$

$$\hat{T} = \text{cov}(X) = \frac{1}{n} X^T X - \Sigma$$

$$\hat{y} = \text{cov}(X, y) = \frac{1}{n} X^T y$$

o When $\Sigma > 0$, \hat{T} is not always positive semidefinite.

o A negative \hat{T} means negative covariance, such as

$$\hat{T} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} - \begin{bmatrix} \sigma^2 I & 0 \\ 0 & \sigma^2 I \end{bmatrix} = \begin{bmatrix} -\sigma^2 I_d & 0 \\ 0 & -\sigma^2 I_d \end{bmatrix}$$

o \hat{T} has negative eigenvalues, $-\sigma^2 I_d$

Ex: o A regularity avoided Saddle point in the eigenvalue space, $\lambda < 0$.

o The solution θ has a singularity at zero.

Example #2: Missing Data

A function category, $\min_{\|\theta\|_1 \leq \sqrt{\frac{n}{\log d}}} \left\{ \frac{1}{2} \theta^T \theta - \langle \bar{y}, \theta \rangle + \lambda_n \|\theta\|_1 \right\}$

$$T = \text{cov}(x) = \frac{X^T X}{n} - p \text{diag}\left(\frac{X^T X}{n}\right)$$

$$\bar{y} = \text{cov}(x, y) = \frac{1}{n} X^T Y$$

- When p shows irregular behavior on the diagonal, specifically $p \neq 0$, but $p \in (0, 1)$.

- T has potential non-convex forms, negative eigenvalues, and bad outcomes for θ .

Example #3: Multiplicative Noise:

A function class, $\min_{\|\theta\|_1 \leq \sqrt{\frac{n}{\log d}}} \left\{ \frac{1}{2} \theta^T T \theta - \langle \bar{y}, \theta \rangle + \lambda_n \|\theta\|_1 \right\}$

$$T = \text{cov}(x) = \frac{1}{n} X^T X / \mathbb{E}[u^T u]$$

$$\bar{y} = \text{cov}(x, y) = \frac{1}{n} X^T Y / \mathbb{E}[u]$$

- u is non-negative noise
- $\mathbb{E}[u^T u]$, covariance in cases generates columns without noise.
- $(\mathbb{E}[u^T u])_{ii}$ columns in a denominator cause singularity, division by zero.
- A case is a multivariate Bernoulli where $\mathbb{E}[u^T u] = \text{cov}(\text{Bernoulli}(p_i, p_j)) = \begin{cases} (1-p_i)(1-p_j) & i \neq j \\ (1-p_i) & i = j \end{cases}$

	Equation	Basic Idea	Examples
Biased Estimator	$E_{X v}[\hat{\theta}] - \theta \neq 0$	<ul style="list-style-type: none"> ◦ Observation ◦ subjective 	<ul style="list-style-type: none"> ◦ Variance ◦ Lower or upper bound ◦ Fehreinheit [Relative]
Unbiased Estimator	$E_{X v}[\hat{\theta}] - \theta = 0$	<ul style="list-style-type: none"> ◦ consistent ◦ objective 	<ul style="list-style-type: none"> ◦ Shrinkage Estimation ◦ Distance ◦ Kelvin [Absolute]

11.12.

a) From Example 11.6, $Z = X + V$, $y = X^T \theta + w$, V is independent of X

$$V = N(x|0, \Sigma_v), \Sigma_v = \text{cov}(v)$$

$$X = N(x|0, \Sigma_x), \Sigma_x = \text{cov}(x)$$

$$\hat{\tau} = \frac{1}{n} Z^T Z - \Sigma_v - \Sigma_x, \hat{\gamma} = \frac{Z^T y}{n}$$

$$\hat{\tau} \text{ and } \hat{\gamma} \text{ are unbiased: } \hat{\tau} - \Sigma_x = \frac{1}{n} Z^T Z - \Sigma_v - \Sigma_x$$

$$\text{Unbiasedness: } = \frac{1}{n} (X + V)^T (X + V) - \Sigma_v - \Sigma_x$$

$$\begin{aligned} &= \cancel{\frac{X^T X}{n}} + 2 \cancel{\frac{X^T V}{n}} + \cancel{\frac{V^T V}{n}} - \cancel{\Sigma_v} - \cancel{\Sigma_x} = 0 \\ &\text{Independent} \end{aligned}$$

$$\hat{\gamma} - \Sigma_{xy} = \frac{Z^T y}{n} - \Sigma_{xy}$$

$$= \frac{(X + V)^T y}{n} - \Sigma_{xy}$$

$$= \frac{X^T y}{n} + \frac{V^T y}{n} - \Sigma_{xy}$$

$$= \frac{X^T y}{n} + \frac{V^T (X\theta + w)}{n} - \Sigma_{xy}$$

$$= \frac{X^T y}{n} - \Sigma_{xy}$$

$$= 0$$

(Proposition 11.10)

$$\|\hat{\tau}\theta^* - \gamma\|_{\max} \leq \phi(Q, \sigma_w) \sqrt{\frac{\log d}{n}}$$

"Covariance subtraction for sub-gaussians becomes subgaussian."

$$11.13. Z_{ij} = \begin{cases} Z_{ij}/1-\nu & \text{if } (i,j) \text{ is observed} \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{\Sigma} = \frac{1}{n} Z^T Z - \nu \text{diag}\left(\frac{Z^T Z}{n}\right); \quad \Sigma_X = \text{cov}(X)$$

a) $\hat{\Sigma}$ and \bar{Y} are unbiased estimators:

$$\begin{aligned} \hat{\Sigma} - \Sigma_X &= \frac{1}{n} Z^T Z - \nu \text{diag}\left(\frac{Z^T Z}{n}\right) - \Sigma_X \\ &= \frac{1}{n} X^T X - \Sigma_X \\ &= 0 \end{aligned}$$

$$\begin{aligned} \bar{Y} - \Sigma_{XY} &= \frac{1}{n} Z^T y - \Sigma_{XY} \\ &= \frac{1}{n} X^T y - \Sigma_{XY} \\ &= 0 \end{aligned}$$

b) X_{ij} and Z_{ij} are sub-Gaussian.

By Proposition 11.13, $P[\|T\theta - Y\|] \leq \underbrace{\phi(Q, \sigma_w)}_{C(\sigma + \|T\theta\|_2)} \sqrt{\frac{\log d}{n}}$

Critical Note: The upper bound in high-dimensional datasets curtail empirical rules in 'normal' distributions.

High-dimensional limit	Empirical Statistics										
$\phi(Q, \sigma_w) \sqrt{\frac{\log d}{n}}$	<table border="1"> <thead> <tr> <th>value</th><th>z-table</th></tr> </thead> <tbody> <tr> <td>one-sigma: σ</td><td>1.0</td></tr> <tr> <td>one: 1.7σ</td><td>1.7</td></tr> <tr> <td>two-sigma: 2σ</td><td>2.0</td></tr> <tr> <td>coefficient: $\sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$</td><td></td></tr> </tbody> </table>	value	z-table	one-sigma: σ	1.0	one: 1.7σ	1.7	two-sigma: 2σ	2.0	coefficient: $\sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$	
value	z-table										
one-sigma: σ	1.0										
one: 1.7σ	1.7										
two-sigma: 2σ	2.0										
coefficient: $\sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$											

b) X_{ij} and V_{ij} are sub-Gaussian

$$\text{By Proposition 11.18, } \mathbb{P}[\|\hat{T}\theta^* - \gamma\|_\infty \leq \underbrace{\phi(Q, \sigma_W)}_{= C(\sigma + \|\theta\|_2)} \sqrt{\frac{\log d}{n}}]$$

$$\mathbb{P}[\|\hat{T}\theta^* - \gamma\|_\infty \leq C(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}]$$

$$\begin{aligned} &= \mathbb{P}[e^{t\hat{T}\theta^* - \gamma}] \stackrel{1 + \lambda \mathbb{E}[\|\hat{T}\theta^* - \gamma\|_\infty] + \frac{\lambda^2 \mathbb{E}[(\|\hat{T}\theta^* - \gamma\|_\infty)^2]}{2}}{e^{\lambda c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}}} \\ &\leq \mathbb{P}[e^{\frac{\lambda^2 \mathbb{E}[(\|\hat{T}\theta^* - \gamma\|_\infty)^2]}{2}}] \leq e^{\lambda c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}} \end{aligned}$$

$$\underset{\lambda}{\operatorname{argmin}} \left\{ \frac{\lambda^2 \mathbb{E}[(\|\hat{T}\theta^* - \gamma\|_\infty)^2]}{2} - \lambda c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}} \right\} = 0$$

$$\lambda^* = \frac{c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}}{\mathbb{E}[(\|\hat{T}\theta^* - \gamma\|_\infty)^2]} = \frac{-c(\sigma + \|\theta\|_2) \frac{\log d}{n}}{2 \mathbb{E}[(\|\hat{T}\theta^* - \gamma\|_\infty)^2]}$$

$$\mathbb{P}[\|\hat{T}\theta^* - \gamma\|_\infty \leq c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}] \leq e^{-\frac{c(\sigma + \|\theta\|_2) \frac{\log d}{n}}{2 \mathbb{E}[(\|\hat{T}\theta^* - \gamma\|_\infty)^2]}}$$

$$\mathbb{E}[\|\hat{T}\theta^* - \gamma\|_\infty] = \frac{1}{\lambda^*} \mathbb{E}[\log e^{-\frac{c(\sigma + \|\theta\|_2) \frac{\log d}{n}}{2 \mathbb{E}[(\|\hat{T}\theta^* - \gamma\|_\infty)^2]}}]$$

$$= \frac{c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}}{2}$$

Note: The (λ^*) goes away with a double-sided dist.

$$c) (\text{Equation 11.50}) \langle \Delta, T\Delta \rangle \geq k \|\Delta\|_2^2 - c_0 \frac{\log d}{n} \|\Delta\|_1^2$$

From definition 7.12, $\frac{1}{n} \|X\Delta\|_2^2 \geq k \|\Delta\|_2^2$ for all Δ

$$\langle \Delta, T\Delta \rangle = \Delta^T \Delta$$

$$\geq k \|\Delta\|_2^2$$

$$\geq k \|\Delta\|_2^2 - c_0 \frac{\log d}{n} \|\Delta\|_2^2$$

$$\begin{aligned} \mathbb{P}\left[\|\hat{T}\theta - \gamma\|_2 \leq c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}\right] \\ := \mathbb{P}\left[e^{1+\lambda\mathbb{E}[\|\hat{T}\theta - \gamma\|_2] + \frac{\lambda^2\mathbb{E}[\|\hat{T}\theta - \gamma\|_2^2]}{2} + \dots} \geq e^{\lambda c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}}\right] \\ \leq \mathbb{P}\left[e^{\frac{\lambda^2\mathbb{E}[\|\hat{T}\theta - \gamma\|_2^2]}{2}} \leq e^{\lambda c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}}\right] \end{aligned}$$

$$\operatorname{argmin}_{\lambda} \left\{ \frac{\lambda^2 \mathbb{E}[\|\hat{T}\theta - \gamma\|_2^2]}{2} - \lambda c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}} \right\} = 0$$

$$\lambda^* = \frac{c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}}{\mathbb{E}[\|\hat{T}\theta - \gamma\|_2^2]}$$

$$\mathbb{P}\left[\|\hat{T}\theta - \gamma\|_2 \leq c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}\right] \leq e^{-\frac{c(\sigma + \|\theta\|_2) \frac{\log d}{n}}{2\mathbb{E}[\|\hat{T}\theta - \gamma\|_2^2]}}$$

$$\mathbb{E}[\|\hat{T}\theta - \gamma\|_2] = \lambda^* \mathbb{E}[\log e^{-\frac{c(\sigma + \|\theta\|_2) \frac{\log d}{n}}{2\mathbb{E}[\|\hat{T}\theta - \gamma\|_2^2]}}]$$

$$= \frac{c(\sigma + \|\theta\|_2) \sqrt{\frac{\log d}{n}}}{2}$$

$$c) (\text{Equation 11.50}) \quad \langle \Delta, T\Delta \rangle \geq k \|\Delta\|_2^2 - C_0 \frac{\log d}{n} \|\Delta\|_1^2$$

From definition 7.12, $\frac{1}{n} \|\Delta\|_2^2 \geq k \|\Delta\|_1^2$ for all Δ

$$\langle \Delta, T\Delta \rangle = \Delta^T \Delta$$

$$\geq k \|\Delta\|_2^2$$

$$\geq k \|\Delta\|_2^2 - C_0 \frac{\log d}{n} \|\Delta\|_1^2$$