

# Linear Regression



Yuling Yan

University of Wisconsin–Madison, Fall 2024

# From classification to regression

---

## Classification:

- there is a joint distribution of  $(X, Y) \sim \rho$  where typically  $X \in \mathbb{R}^d$  and  $Y \in \{1, \dots, K\}$  is discrete
- Goal: given input  $x$ , find the label  $y$  with the highest posterior probability

$$\arg \max_{y \in \{1, \dots, K\}} \mathbb{P}(Y = y | X = x)$$

## Regression:

- there is a joint distribution of  $(X, Y) \sim \rho$  where  $X \in \mathbb{R}^d$  and  $Y \in \mathbb{R}$
- Goal: given input  $x$ , find a prediction  $f(x)$  for  $Y$  conditional on  $X = x$ , that minimizes MSE

$$\mathbb{E}[(Y - f(x))^2 | X = x]$$

# Target of regression problem

---

## Theorem 4.1

For any random variable  $Z$ , we have

$$\arg \min_{c \in \mathbb{R}} \mathbb{E}[(Z - c)^2] = \mathbb{E}[Z].$$

### Implications for regression problem:

- Conditional on  $X = x$ , the optimal prediction for  $Y$  that minimizes MSE is

$$f^*(x) = \mathbb{E}[Y|X = x]$$

- Rewrite the model

$$Y = \underbrace{\mathbb{E}[Y|X]}_{\text{regression function}} + \underbrace{Y - \mathbb{E}[Y|X]}_{\text{mean-zero noise}}$$

# Regression problem

---

We will consider the regression problem in a more straightforward way:

$$y = f^*(\mathbf{x}) + \varepsilon$$

- $\mathbf{x} \in \mathbb{R}^d$  is the input,  $y \in \mathbb{R}$  is the output
- $\varepsilon$  is some mean-zero random noise, e.g.,  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$
- $f^* : \mathbb{R}^d \rightarrow \mathbb{R}$  is the *unknown* regression function
- Training data:  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  satisfying

$$y_i = f^*(\mathbf{x}_i) + \varepsilon_i$$

where  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d. noise with  $\mathbb{E}[\varepsilon_i] = 0$ , and

- in some cases, we assume  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are deterministic (**fixed design**)
  - sometimes we may assume that  $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{i.i.d.}}{\sim} \rho_X$  (**random design**)
- Learn the regression function  $f^*$  based on training data

# Overview

---

- **Linear regression:** model the regression function  $f^*$  as a linear function

$$f^*(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\beta}^*$$

where we assume  $\mathbf{x}$  includes a constant variable 1. Here  $\boldsymbol{\beta}^* \in \mathbb{R}^d$  is the unknown parameter.

- **Nonparametric regression:** assume that

$$f^* \in \mathcal{F}$$

where  $\mathcal{F}$  is certain function class, e.g.,

- class of quadratic function
- class of convex function
- Reproducing Kernel Hilbert Space (RKHS)

*Linear regression: classical setting*

# Linear regression

---

- Linear regression:

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta}^* + \varepsilon_i \quad (i = 1, \dots, n)$$

where  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are fixed design, and  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d. noise satisfying  $\mathbb{E}[\varepsilon_i] = 0$  and  $\text{var}(\varepsilon_i) = \sigma^2$

- Consider matrix notation

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\varepsilon}$$

where

$$\mathbf{Y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix} \in \mathbb{R}^{n \times d}, \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} \in \mathbb{R}^n$$

# Least square estimator

---

- The most popular estimation method is *least squares*, which estimates  $\beta^*$  by minimizing the residual sum of squares

$$\sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2.$$

- Ordinary least squares (OLS) estimator:

$$\hat{\boldsymbol{\beta}} := \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^d} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$$

It has minimizer

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}.$$

- Suppose the noise are i.i.d. Gaussian, then OLS is the MLE



# Theoretical properties

- Linear estimator: estimator of the form  $\mathbf{A}\mathbf{Y}$  for some matrix  $\mathbf{A} \in \mathbb{R}^{d \times n}$
- OLS achieves the minimum variance among all linear unbiased estimators
- Furthermore, when the noise is i.i.d. Gaussian, OLS achieves the minimum variance among all unbiased estimators

## Theorem 4.2

- **Gauss-Markov:** *The OLS estimator  $\hat{\beta}$  is the best linear unbiased estimator of  $\beta^*$ , i.e. for any linear and unbiased estimator  $\tilde{\beta}$  of  $\beta^*$ ,*

$$\text{cov}(\hat{\beta}) \preceq \text{cov}(\tilde{\beta}).$$

- **Cramér-Rao lower bound:** *when  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d.  $\mathcal{N}(0, \sigma^2)$ , the variance of OLS matches the Cramér-Rao lower bound, i.e. for any unbiased estimator  $\tilde{\beta}$  of  $\beta^*$ ,*

$$\text{cov}(\hat{\beta}) \preceq \text{cov}(\tilde{\beta}).$$

# Cramér-Rao lower bound

---

- Consider  $X_1, \dots, X_n$  be i.i.d. samples from a density  $f_\theta$
- The unknown parameter  $\theta \in \Theta$
- Let  $T(X_1, \dots, X_n)$  be any unbiased estimator for  $\theta$
- Under some regularity condition,

$$\text{cov}(T(X_1, \dots, X_n)) \succeq [I(\theta)]^{-1}$$

where  $I(\theta)$  is the **Fisher information matrix**

$$\begin{aligned} I(\theta) &= n\mathbb{E}_{X \sim f_\theta} [\nabla_\theta \log f_\theta(X) [\nabla_\theta \log f_\theta(X)]^\top] \\ &= -n\mathbb{E}_{X \sim f_\theta} [\nabla_\theta^2 \log f_\theta(X)] \end{aligned}$$

# Implications

---

- The OLS estimator is the best one among all unbiased estimator for  $\beta^*$  in terms of minimizing MSE (why?)
- Is it also the best estimator among any estimator for  $\beta^*$ , including those biased ones?

# Implications

---

- The OLS estimator is the best one among all unbiased estimator for  $\beta^*$  in terms of minimizing MSE (why?)
- Is it also the best estimator among any estimator for  $\beta^*$ , including those biased ones?
  - *No! There are biased estimator which can achieve smaller MSE.*

# Implications

---

- The OLS estimator is the best one among all unbiased estimator for  $\beta^*$  in terms of minimizing MSE (why?)
- Is it also the best estimator among any estimator for  $\beta^*$ , including those biased ones?
  - *No! There are biased estimator which can achieve smaller MSE.*
- Examples of biased estimator with smaller MSE:
  - James-Stein estimator
  - Ridge regression

— *shrinkage estimators*

*Shrinkage estimator*

# Bias-variance tradeoff

---

- Suppose that the unknown parameter is  $\beta^* \in \mathbb{R}^d$
- For any estimator  $\hat{\beta}$  (more generally, any random vector), the mean squared error (MSE) can be decomposed into

$$\underbrace{\mathbb{E}[\|\hat{\beta} - \beta^*\|_2^2]}_{=: \text{MSE}} = \underbrace{\|\mathbb{E}[\hat{\beta}] - \beta^*\|_2^2}_{\text{bias}} + \underbrace{\text{tr}(\text{cov}(\hat{\beta}))}_{\text{variance}}$$

- For unbiased estimator (e.g., OLS), the bias is zero
- By tolerating a small amount of bias we may be able to achieve a larger reduction in variance, thus achieving smaller MSE

# James-Stein estimator

---

- Consider a Gaussian sequence model,

$$\mathbf{Y} = \boldsymbol{\beta}^* + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$$

which is a special linear regression by taking  $d = n$  and  $\mathbf{X} = \mathbf{I}_n$

- OLS / MLE:  $\hat{\boldsymbol{\beta}}_{\text{OLS}} = \mathbf{Y}$
- James-Stein estimator:

$$\hat{\boldsymbol{\beta}}_{\text{JS}} = \left(1 - \frac{n-2}{\|\mathbf{Y}\|_2^2}\right) \mathbf{Y}$$

## Theorem 4.3

*James-Stein estimator has smaller MSE than OLS when  $n \geq 3$ , i.e.,*

$$\text{MSE}(\hat{\boldsymbol{\beta}}_{\text{JS}}) < \text{MSE}(\hat{\boldsymbol{\beta}}_{\text{OLS}}) \quad \text{for any } \boldsymbol{\beta}^*$$



# Implications

---

- By shrinking the OLS towards zero, we achieve smaller MSE  
— *inadmissibility of OLS (or MLE)*
- It is not even necessary to shrink towards zero: for any fixed  $\mathbf{c} \in \mathbb{R}^n$ ,

$$\hat{\beta}_{\mathbf{c}} := \mathbf{Y} - \frac{p-2}{\|\mathbf{Y} - \mathbf{c}\|_2^2} (\mathbf{Y} - \mathbf{c})$$

also satisfy the same property as Theorem 4.3

- Can be extended to linear regression:

$$\hat{\beta}_{\text{JS}} = \hat{\beta}_{\text{OLS}} - \frac{(d-2)\hat{\sigma}^2}{\|\mathbf{X}^\top \mathbf{X} \hat{\beta}_{\text{OLS}}\|_2^2} \mathbf{X}^\top \mathbf{X} \hat{\beta}_{\text{OLS}}.$$

# Ridge regression

---

- Ridge regression:  $\ell_2$ -penalized least squares estimator

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2,$$

where  $\lambda$  is the tuning parameter.

- The ridge regression estimator admits closed-form solution:

$$\hat{\beta}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{Y}.$$

It is well defined even when  $\mathbf{X}^\top \mathbf{X}$  is not invertible

- As  $\lambda \rightarrow 0$ , ridge regression recovers the OLS
- Interpretation as MAP estimator with a Gaussian prior on  $\beta^*$

# MAP estimate

---

Consider observing  $X$  from a density  $f_{\theta^*}$ , where  $\theta^* \in \Theta$  is unknown

**Frequentist's viewpoint:**  $\theta^*$  is fixed (though unknown)

- Likelihood function:  $f_{\theta}(X)$  (a function of  $\theta \in \Theta$ )
- Estimate  $\theta^*$  by the maximizer of the likelihood function
  - *maximum likelihood estimation (MLE)*

**Bayesian's viewpoint:**  $\theta$  is also random

- We have a prior distribution  $g(\theta)$  over  $\Theta$ , and conditional on  $\theta$ ,  $X \sim f_{\theta}$
- Posterior probability of  $\theta$  after observing  $X$ :

$$\mathbb{P}(\theta|X) = \frac{g(\theta)f_{\theta}(X)}{\int_{\Theta} g(\theta')f_{\theta'}(X)d\theta'} \propto g(\theta)f_{\theta}(X)$$

- Estimate  $\theta$  by the maximizer of the posterior probability
  - *maximum a posteriori estimation (MAP)*

# Properties of ridge regression

---

**Ridge regression:**

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2 = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{Y}.$$

## Theorem 4.4

*There exists  $\lambda_0 > 0$  such that ridge regression  $\hat{\beta}_\lambda$  achieves smaller MSE than OLS estimate*

$$\text{MSE}(\hat{\beta}_\lambda) < \text{MSE}(\hat{\beta}_{\text{OLS}})$$

*for any  $\lambda \in (0, \lambda_0]$ .*

# Properties of ridge regression

---

**Ridge regression:**

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^d} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2 = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{Y}.$$

## Theorem 4.4

*There exists  $\lambda_0 > 0$  such that ridge regression  $\hat{\beta}_\lambda$  achieves smaller MSE than OLS estimate*

$$\text{MSE}(\hat{\beta}_\lambda) < \text{MSE}(\hat{\beta}_{\text{OLS}})$$

*for any  $\lambda \in (0, \lambda_0]$ .*

- To prove this theorem, we need some tool from linear algebra

# Singular Value Decomposition (SVD)

---

For any rank- $r$  matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , it can be expressed as

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$$

- $\mathbf{U} \in \mathbb{R}^{n \times r}$  and  $\mathbf{V} \in \mathbb{R}^{d \times r}$  are orthogonal matrices:

$$\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_r], \quad \mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_r],$$

where  $\{\mathbf{u}_i\}_{i=1}^r$  (resp.  $\{\mathbf{v}_i\}_{i=1}^r$ ) are orthonormal vectors in  $\mathbb{R}^n$  (resp.  $\mathbb{R}^d$ )

- $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$  is a diagonal matrix

$$\mathbf{\Sigma} = \text{diag}\{\sigma_1, \dots, \sigma_r\}$$

where  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  are the singular values of  $\mathbf{X}$

## More about SVD

---

For any rank- $r$  matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  with SVD  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$

- Connection to eigen-decomposition

$$\mathbf{X}\mathbf{X}^\top = \mathbf{U}\mathbf{\Sigma}^2\mathbf{U}^\top = \begin{bmatrix} \mathbf{U} & \mathbf{U}_\perp \end{bmatrix} \begin{bmatrix} \mathbf{\Sigma}^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{n-r} \end{bmatrix} \begin{bmatrix} \mathbf{U}^\top \\ \mathbf{U}_\perp^\top \end{bmatrix}$$

$$\mathbf{X}^\top\mathbf{X} = \mathbf{V}\mathbf{\Sigma}^2\mathbf{V}^\top = \begin{bmatrix} \mathbf{V} & \mathbf{V}_\perp \end{bmatrix} \begin{bmatrix} \mathbf{\Sigma}^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{d-r} \end{bmatrix} \begin{bmatrix} \mathbf{V}^\top \\ \mathbf{V}_\perp^\top \end{bmatrix}$$

where  $\mathbf{U}_\perp$  (resp.  $\mathbf{V}_\perp$ ) is the orthogonal complement of  $\mathbf{U}$  (resp.  $\mathbf{V}$ )

- The operator (spectral) norm of  $\mathbf{X}$

$$\|\mathbf{X}\| = \sup_{\|\mathbf{a}\|_2=1} \|\mathbf{X}\mathbf{a}\|_2 = \sigma_1$$

- The Frobenius norm of  $\mathbf{X}$

$$\|\mathbf{X}\|_F^2 = \sum_{i=1}^r \sigma_i^2$$

# Implications to ridge regression

---

Suppose that the design matrix  $X$  has SVD  $U\Sigma V^\top$

- Bias-variance decomposition

$$\mathbb{E}[\|\hat{\beta}_\lambda - \beta^*\|_2^2] = \|\mathbb{E}[\hat{\beta}_\lambda] - \beta^*\|_2^2 + \text{tr}(\text{cov}(\hat{\beta}_\lambda))$$

- Bias term

$$\|\mathbb{E}[\hat{\beta}_\lambda] - \beta^*\|_2^2 = \sum_{i=1}^d \left( \frac{\lambda \tilde{\beta}_i}{\lambda + \sigma_i^2} \right)^2 \quad \text{where } \tilde{\beta} = [V, V_\perp]^\top \beta^*$$

- Variance term

$$\text{cov}(\hat{\beta}_\lambda) = \sigma^2 \sum_{i=1}^d \left( \frac{\sigma_i}{\lambda + \sigma_i^2} \right)^2$$

- This allows us to prove Theorem 4.4



*Linear regression: high-dimensional setting*

# What happens in high-dimension?

---

**High-dimensional linear regression:**

$$Y = X\beta^* + \varepsilon$$

where the dimension  $d$  is much larger than the sample size  $n$

- OLS fails because  $X^T X$  is not invertible
- In general, it is not possible to say something meaningful about  $\beta^* \in \mathbb{R}^d$  from  $n$  samples  $Y \in \mathbb{R}^n$  (identifiability issue)

# What happens in high-dimension?

---

## High-dimensional linear regression:

$$Y = X\beta^* + \varepsilon$$

where the dimension  $d$  is much larger than the sample size  $n$

- OLS fails because  $X^\top X$  is not invertible
- In general, it is not possible to say something meaningful about  $\beta^* \in \mathbb{R}^d$  from  $n$  samples  $Y \in \mathbb{R}^n$  (identifiability issue)
- A meaningful and workable setup: assume  $\beta^*$  is sparse, i.e.,

$$s := \|\beta^*\|_0 \equiv |\{j : \beta_j^* \neq 0\}| \ll d$$

# Sparse linear regression

---

**High-dimensional linear regression:**

$$Y = X\beta^* + \epsilon$$

where  $d \geq n$ , but  $s = \|\beta^*\|_0 \ll d$

- **Genomics:** only a small subset of genes is expected to be associated with a particular trait or disease
- **Finance and Economics:** only a small subset of macroeconomic variables or market signals may be relevant to stock returns or economic growth
- .....

# Insights

---

- Motivated by ridge regression, we may consider

$$\arg \min_{\beta \in \mathbb{R}^d} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_0$$

- Issue: computationally hard ( $\|\cdot\|_0$  is discontinuous, non-convex...)
- Idea: use  $\|\cdot\|_1$  instead
- Insights from compressed sensing (noiseless): under certain conditions (known as restricted isometry property),  $\ell_1$  minimization problem

$$\arg \min_{\beta \in \mathbb{R}^d} \|\beta\|_1 \quad \text{s.t.} \quad \mathbf{X}\beta = \mathbf{Y}$$

has unique minimizer that coincides with the minimizer to

$$\arg \min_{\beta \in \mathbb{R}^d} \|\beta\|_0 \quad \text{s.t.} \quad \mathbf{X}\beta = \mathbf{Y}.$$

# LASSO

---

LASSO (Least Absolute Shrinkage and Selection Operator) estimates  $\beta^*$  by solving the following convex optimization problem:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^d} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1,$$

where:

- $\|\mathbf{Y} - \mathbf{X}\beta\|_2^2$ : residual sum of squares (RSS).
- $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$ :  $\ell_1$ -norm penalty.
- $\lambda > 0$ : tuning parameter that controls the trade-off between **goodness of fit** and **sparsity**.
- Interpretation as MAP estimator with a Laplace prior on  $\beta^*$
- Questions:
  - How to compute LASSO estimate?
  - What is the statistical properties of LASSO?

*How to compute LASSO: proximal gradient method*

# A more general class of convex optimization

---

Consider unconstrained convex optimization problem of the form

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) := f(\mathbf{x}) + h(\mathbf{x})$$

where

- $f(\mathbf{x})$ : a differentiable, convex function
- $h(\mathbf{x})$ : a convex, potentially non-differentiable function (e.g.,  $\ell_1$ -norm).
- Example: LASSO can be viewed as taking

$$f(\mathbf{x}) = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2, \quad h(\mathbf{x}) = \lambda \|\boldsymbol{\beta}\|_1.$$

**Issue:** gradient descent (GD) does not work (due to non-smoothness)



# A Proximal View of Gradient Descent

---

- To motivate proximal gradient methods, we first revisit gradient descent for  $\min_{\mathbf{x}} f(\mathbf{x})$ , where  $f(\cdot)$  is convex and smooth
- Gradient descent update:  $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$
- This is equivalent to

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x}} \left\{ \underbrace{f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle}_{\text{first-order approximation at } \mathbf{x}_t} + \underbrace{\frac{1}{2\eta} \|\mathbf{x} - \mathbf{x}_t\|_2^2}_{\text{proximal term}} \right\}$$

- Heuristics: search for  $\mathbf{x}_{t+1}$  that
  - aim to minimize  $f(\cdot)$  (through minimizing first-order approximation)
  - remains close to  $\mathbf{x}_t$  such that first-order approximation at  $\mathbf{x}_t$  is valid (enforced by proximal term)
- Benefit: minimizing a quadratic function, admits simple solution (i.e., GD)

# Proximal gradient method: algorithm

---

Consider an iterative algorithm: starting from  $\mathbf{x}_t$ , update

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x}} \left\{ \underbrace{f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle}_{\text{first-order approximation at } \mathbf{x}_t} + h(\mathbf{x}) + \underbrace{\frac{1}{2\eta} \|\mathbf{x} - \mathbf{x}_t\|_2^2}_{\text{proximal term}} \right\}$$

- Define proximal operator

$$\text{prox}_h(\mathbf{v}) = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ h(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|_2^2 \right\}$$

- If this proximal operator is easy to compute, then we can express

$$\mathbf{x}_{t+1} = \text{prox}_{\eta h}(\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t))$$

- alternates between gradient updates on  $f$  and proximal minimization on  $h$

# Proximal gradient method: properties

---

**Proximal gradient algorithm:** for  $t = 1, 2, \dots$

$$\mathbf{x}_{t+1} = \text{prox}_{\eta h}(\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t))$$

- fast convergence when  $f$  is convex and  $L$ -smooth: take  $\eta = 1/L$ ,

$$F(\mathbf{x}_t) - F^* \leq \frac{L}{2t} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$$

- exponential convergence when  $f$  is  $\mu$ -strongly convex

$$\|\mathbf{x}_t - \mathbf{x}^*\|_2^2 \leq (1 - \mu/L)^t \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$$

- when  $h(\mathbf{x}) = 0$  when  $\mathbf{x} \in \mathcal{A}$  and  $h(\mathbf{x}) = \infty$  otherwise, this gives the projected gradient descent for  $\min_{\mathbf{x} \in \mathcal{A}} f(\mathbf{x})$ :

$$\mathbf{x}_{t+1} = \mathcal{P}_{\mathcal{A}}(\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t))$$

- Recommended reading material: Lecture 5 of the course [Large-Scale Optimization for Data Science](#)

# Application to LASSO

---

- LASSO:

$$f(\boldsymbol{\beta}) = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \quad \text{and} \quad h(\boldsymbol{\beta}) = \lambda\|\boldsymbol{\beta}\|_1$$

- The proximal operator admits closed-form expression

$$\text{prox}_h(\mathbf{v}) = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\boldsymbol{\beta} - \mathbf{v}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\} = \text{shrink}_\lambda(\mathbf{v})$$

where  $\text{shrink}_\lambda(\cdot)$  applies entrywise shrinkage to  $\mathbf{v}$  towards zero:

$$[\text{shrink}_\lambda(\mathbf{v})]_j = \begin{cases} v_j - \lambda, & \text{if } v_j \geq \lambda, \\ v_j + \lambda, & \text{if } v_j \leq -\lambda, \\ 0, & \text{otherwise.} \end{cases}$$

- Proximal gradient algorithm for LASSO:

$$\boldsymbol{\beta}_{t+1} = \text{shrink}_{\eta\lambda}(\boldsymbol{\beta}_t - 2\eta\mathbf{X}^\top\mathbf{X}\boldsymbol{\beta}_t + 2\eta\mathbf{X}^\top\mathbf{Y})$$

## *Statistical properties of LASSO*

# Setup

---

## LASSO:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 \right\},$$

- Independent, sub-Gaussian noise  $\|\varepsilon_i\|_{\psi_2} \leq \sigma$
- Sparsity:  $n \gg s \log d$
- Theory-informed tuning parameter selection:

$$\lambda \asymp \sigma \sqrt{n \log d}$$

- Question:
  - Does LASSO recover the support of  $\beta^*$ ?
  - Does LASSO provide reliable estimate for  $\beta^*$ ?

# Optimality condition

---

The optimality condition for unconstrained convex optimization

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$$

- if  $f$  is smooth:  $\nabla f(\hat{\mathbf{x}}) = \mathbf{0}$
- in general (when  $f$  might not be smooth):  $\mathbf{0} \in \partial f(\hat{\mathbf{x}})$

Here  $\partial f(\mathbf{x}) \subseteq \mathbb{R}^d$  is the **subgradient** of the convex function  $f$  at  $\mathbf{x}$ :

$$\mathbf{g} \in \partial f(\mathbf{x}) \iff f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x}) \quad \text{for all } \mathbf{y} \in \mathbb{R}^d$$

Check (in homework):

- if  $f$  is smooth at  $\mathbf{x}$ :  $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$
- the optimality condition for LASSO is: for each  $1 \leq j \leq d$

$$[\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}^\top \hat{\boldsymbol{\beta}})]_j \quad \begin{cases} = \lambda \cdot \text{sign}(\hat{\beta}_j) & \text{if } \hat{\beta}_j \neq 0 \\ \in [-\lambda, \lambda] & \text{if } \hat{\beta}_j = 0 \end{cases}$$

# Model selection consistency

---

- Let  $S = \{j : \beta_j^* \neq 0\}$  be the support set (nonzero coefficients) and  $S^c$  be its complement.
- **Irrepresentable condition:**

$$\|\mathbf{X}_{S^c}^\top \mathbf{X}_S (\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \boldsymbol{\beta}_S^*\|_\infty < 1,$$

where  $\mathbf{X}_S$  and  $\mathbf{X}_{S^c}$  as submatrices of  $\mathbf{X}$  with columns corresponding to  $S$  and  $S^c$ , and  $\boldsymbol{\beta}_S^*$  is the sub-vector of  $\boldsymbol{\beta}^*$  corresponding to  $S$

- **Model Selection Consistency:** If the irrepresentable condition holds, under certain assumptions, the Lasso estimator satisfies:

$$\mathbb{P}(\widehat{S} = S) \rightarrow 1 \quad \text{as } n \rightarrow \infty,$$

where  $\widehat{S} = \{j : \widehat{\beta}_j \neq 0\}$ .



# Estimation guarantees

---

- **Restricted eigenvalue condition:** For any  $\mathbf{v} \in \mathbb{R}^p$  such that  $\|\mathbf{v}_{S^c}\|_1 \leq 3\|\mathbf{v}_S\|_1$ , the restricted eigenvalue condition is:

$$\min_{\|\mathbf{v}\|_2=1, \|\mathbf{v}_{S^c}\|_1 \leq 3\|\mathbf{v}_S\|_1} \mathbf{v}^\top \left( \frac{1}{n} \mathbf{X}^\top \mathbf{X} \right) \mathbf{v} > 0.$$

This is satisfied by e.g., i.i.d. Gaussian matrix  $\mathbf{X}$ .

- **Estimation error:** If the restricted eigenvalue condition holds, under certain assumptions, the LASSO estimator satisfies:

$$\frac{1}{n} \|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)\|_2^2 \lesssim \sigma^2 s \frac{\log d}{n},$$

and

$$\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_1 \lesssim \sigma s \sqrt{\frac{\log d}{n}}.$$

# Reference

---

## Model selection:

- Peng Zhao, and Bin Yu. "On model selection consistency of Lasso." *The Journal of Machine Learning Research* 7 (2006): 2541-2563.
- Martin J. Wainwright. "Sharp thresholds for High-Dimensional and noisy sparsity recovery using  $\ell_1$ -Constrained Quadratic Programming (Lasso)." *IEEE transactions on information theory* 55.5 (2009): 2183-2202.

## Estimation error bounds:

- Peter J. Bickel, Ya'acov Ritov, and Alexandre B. Tsybakov. "Simultaneous analysis of Lasso and Dantzig selector." *Annals of Statistics* 37.4 (2009): 1705-1732.