

Regularization

Linear regression when $d > n$

Data in matrix/vector form

$$\mathbf{A} := \underbrace{\begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ \leftarrow & \mathbf{x}_2^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix}}_{n \times d \text{ matrix}}, \quad \mathbf{b} := \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}}_{n \times 1 \text{ vector}}.$$

Ordinary least squares

Ordinary least squares $\hat{\mathbf{w}}_{\text{ols}}$: typically “defined” by $\hat{\mathbf{w}}_{\text{ols}} := (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$.
Ill-defined when $d > n$. In this case, the $d \times d$ matrix

$$\mathbf{A}^\top \mathbf{A} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$$

is not invertible (as its rank is at most $n < d$).

1 / 29

2 / 29

Regularization

Understanding regularization

Typical solution: regularization

Some examples:

- ▶ encourage $\|\mathbf{w}\|_2^2$ to be small (“ridge regression”)
- ▶ encourage $\|\mathbf{w}\|_1$ to be small (“Lasso”)
- ▶ encourage \mathbf{w} to be sparse (“sparse regression”)

Example: regularized least squares criterion

Find $\mathbf{w} \in \mathbb{R}^d$ to minimize

$$\|\mathbf{A}\mathbf{w} - \mathbf{b}\|_2^2 + \lambda R(\mathbf{w})$$

where $\lambda > 0$ and $R: \mathbb{R}^d \rightarrow \mathbb{R}_+$ is a *penalty function* / *regularizer*.

How do I decide which type of regularization to use?

- ▶ Could just try them all ...
- ▶ Better answer: try to understand their statistical behavior in a broad class of scenarios.

3 / 29

4 / 29

Ridge regression

Ridge regression

Ridge regression

Find $\mathbf{w} \in \mathbb{R}^d$ to minimize

$$\|\mathbf{A}\mathbf{w} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$.

This *always* has a unique solution.

6 / 29

Ridge regression via calculus

Ridge regression objective is convex function of \mathbf{w} .
Suffices to find \mathbf{w} where gradient is zero.

$$\nabla_{\mathbf{w}} \left\{ \|\mathbf{A}\mathbf{w} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \right\} = 2\mathbf{A}^\top (\mathbf{A}\mathbf{w} - \mathbf{b}) + 2\lambda \mathbf{w}.$$

This is zero when

$$(\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I}) \mathbf{w} = \mathbf{A}^\top \mathbf{b},$$

a system of linear equations in \mathbf{w} .

Matrix $\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I}$ is invertible since $\lambda > 0$, so its *unique* solution is

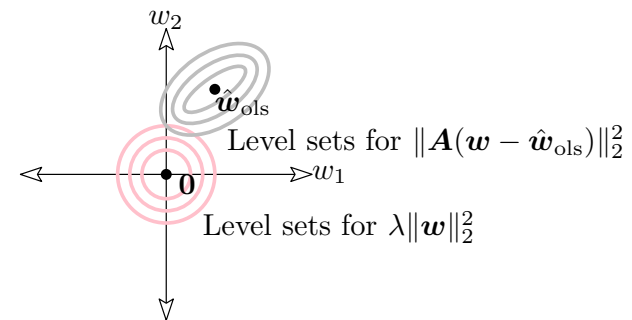
$$\hat{\mathbf{w}}_{\text{ridge}} := (\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^\top \mathbf{b}.$$

7 / 29

Ridge regression: geometry

If $\hat{\mathbf{w}}_{\text{ols}}$ exists, then ridge regression objective (as function of \mathbf{w}) is

$$\|\mathbf{A}(\mathbf{w} - \hat{\mathbf{w}}_{\text{ols}})\|_2^2 + \lambda \|\mathbf{w}\|_2^2 + (\text{stuff not depending on } \mathbf{w}).$$



8 / 29

Aside: Eigendecompositions

Every symmetric matrix $M \in \mathbb{R}^{d \times d}$ guaranteed to have eigendecomposition with real eigenvalues:

M

$(d \times d)$

=

V

$(d \times d)$

Λ

$(d \times d)$

V^\top

$(d \times d)$

=

$\sum_{i=1}^d \lambda_i v_i v_i^\top$

real **eigenvalues**: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ ($\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$);
corresponding orthonormal **eigenvectors**: v_1, v_2, \dots, v_d ($V = [v_1 | v_2 | \dots | v_d]$).

Eigenvectors v_1, v_2, \dots, v_d constitute an **orthonormal basis** for \mathbb{R}^d .

So every $w \in \mathbb{R}^d$ can be written as a linear combination of these vectors:

$$w = \sum_{j=1}^d \langle v_j, w \rangle v_j.$$

9 / 29

Ridge regression: eigendecomposition

Write eigendecomposition of $A^\top A$ as

$$A^\top A = \sum_{j=1}^d \lambda_j v_j v_j^\top$$

where $v_1, v_2, \dots, v_d \in \mathbb{R}^d$ are orthonormal eigenvectors with corresponding eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$.

► For $\lambda > 0$, the inverse of $A^\top A + \lambda I$ exists, and has the form

$$(A^\top A + \lambda I)^{-1} = \sum_{j=1}^d \frac{1}{\lambda_j + \lambda} v_j v_j^\top.$$

10 / 29

Ridge regression vs. ordinary least squares

If \hat{w}_{ols} exists, then

$$\begin{aligned} \hat{w}_{ridge} &= (A^\top A + \lambda I)^{-1} (A^\top A) \hat{w}_{ols} \\ &= \left(\sum_{j=1}^d \frac{1}{\lambda_j + \lambda} v_j v_j^\top \right) \left(\sum_{j=1}^d \lambda_j v_j v_j^\top \right) \hat{w}_{ols} \\ &= \left(\sum_{j=1}^d \frac{\lambda_j}{\lambda_j + \lambda} v_j v_j^\top \right) \hat{w}_{ols} \quad (\text{by orthogonality}) \\ &= \sum_{j=1}^d \frac{\lambda_j}{\lambda_j + \lambda} \langle v_j, \hat{w}_{ols} \rangle v_j. \end{aligned}$$

Interpretation: Shrink \hat{w}_{ols} towards zero by $\frac{\lambda_j}{\lambda_j + \lambda}$ factor in direction v_j .

Effective degrees-of-freedom: $df(\lambda) := \sum_{j=1}^d \frac{\lambda_j}{\lambda_j + \lambda}$.

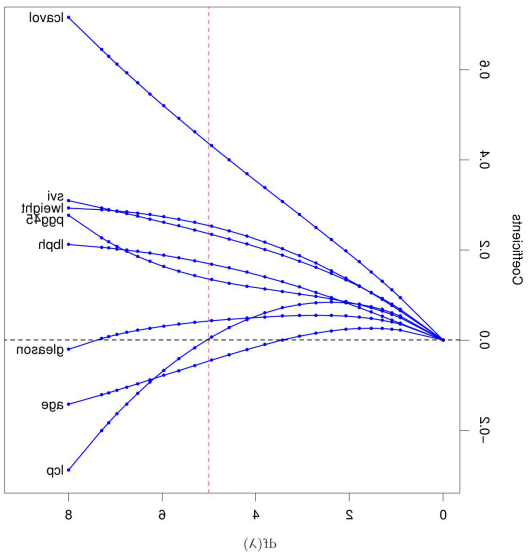
11 / 29

Coefficient profile

Horizontal axis: varying λ (large λ to left, small λ to right).
Vertical axis: coefficient value in \hat{w}_{ridge} for eight different variables.

12 / 29

Coefficient profile (flipped)



Horizontal axis: varying λ (small λ to left, large λ to right).
Vertical axis: coefficient value in $\hat{\mathbf{w}}_{\text{ridge}}$ for eight different variables.

How to compare OLS and ridge?

- ▶ **Case 1:** $\mathbf{A}^\top \mathbf{A}$ not invertible.
No comparison, because OLS even doesn't exist.
- ▶ **Case 2:** $\mathbf{A}^\top \mathbf{A}$ is invertible (but perhaps close to being singular).
How do ridge regression and OLS compare?

Is there a general setting in which to analyze OLS and ridge regression?
 - ▶ **Statistical learning setting:** data comes from unknown distribution P over $\mathbb{R}^d \times \mathbb{R}$, goal is to minimize expected square loss.
(In context of regression, typically called **random design**.)
 - ▶ ...

Fixed-design setting

Easier/cleaner to study OLS and ridge regression in the **fixed design** setting:

- ▶ Assume $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are not random; only y_1, y_2, \dots, y_n are random.
(\mathbf{A} is not random; only \mathbf{b} is random.)

Also assume all y_i have finite variance.

- ▶ Best predictor of y_i is $\mathbb{E}[y_i]$ (in terms of expected square loss).
Note: $\mathbf{x}_i \mapsto \mathbb{E}[y_i]$ might not be realized by a linear function.
- ▶ Let $\mathbf{w}_* \in \mathbb{R}^d$ be a weight vector that minimizes

$$\mathbf{w} \mapsto \mathbb{E}[\|\mathbf{A}\mathbf{w} - \mathbf{b}\|_2^2]$$

(best *linear* predictor).

If $\mathbf{A}^\top \mathbf{A}$ is invertible, then $\mathbf{w}_* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbb{E}[\mathbf{b}]$.

- ▶ **Basic question:** how well can we estimate \mathbf{w}_* ?

OLS: fixed-design analysis

When $\hat{\mathbf{w}}_{\text{ols}}$ exists (i.e., when $\mathbf{A}^\top \mathbf{A}$ is invertible), then it is an unbiased estimator of \mathbf{w}_* :

$$\mathbb{E}[\hat{\mathbf{w}}_{\text{ols}}] = \mathbf{w}_*.$$

Let $\boldsymbol{\varepsilon} := \mathbf{b} - \mathbf{A}\mathbf{w}_*$. Then

$$\text{cov}(\hat{\mathbf{w}}_{\text{ols}}) = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbb{E}[\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top] \mathbf{A} (\mathbf{A}^\top \mathbf{A})^{-1}.$$

For example, if $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, then

$$\text{cov}(\hat{\mathbf{w}}_{\text{ols}}) = \sigma^2 (\mathbf{A}^\top \mathbf{A})^{-1} = \sigma^2 \sum_{j=1}^d \frac{1}{\lambda_j} \mathbf{v}_j \mathbf{v}_j^\top,$$

so the variance of $\langle \mathbf{v}_j, \hat{\mathbf{w}}_{\text{ols}} \rangle$ is

$$\text{var}(\langle \mathbf{v}_j, \hat{\mathbf{w}}_{\text{ols}} \rangle) = \mathbf{v}_j^\top \text{cov}(\hat{\mathbf{w}}_{\text{ols}}) \mathbf{v}_j = \frac{\sigma^2}{\lambda_j}.$$

Note: if λ_d is very close to zero (so $\mathbf{A}^\top \mathbf{A}$ is close to being singular), then the variance in direction \mathbf{v}_d is very high.

Ridge regression: fixed-design analysis

Ridge regression is not an *unbiased* estimator \mathbf{w}_* :

$$\mathbb{E}[\hat{\mathbf{w}}_{\text{ridge}}] \neq \mathbf{w}_*.$$

But, covariance of $\hat{\mathbf{w}}_{\text{ridge}}$ is always “smaller” than that of $\hat{\mathbf{w}}_{\text{ols}}$:

$$\begin{aligned} \text{cov}(\hat{\mathbf{w}}_{\text{ridge}}) &= (\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^\top \mathbb{E}[\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top] \mathbf{A} (\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I})^{-1} \\ &= \mathbf{M} \text{cov}(\hat{\mathbf{w}}_{\text{ols}}) \mathbf{M} \quad (\text{if } \hat{\mathbf{w}}_{\text{ols}} \text{ exists}), \end{aligned}$$

where

$$\mathbf{M} := \sum_{j=1}^d \frac{\lambda_j}{\lambda_j + \lambda} \mathbf{v}_j \mathbf{v}_j^\top.$$

For example, if $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, then

$$\text{cov}(\hat{\mathbf{w}}_{\text{ridge}}) = \sigma^2 \sum_{j=1}^d \frac{\lambda_j}{(\lambda_j + \lambda)^2} \mathbf{v}_j \mathbf{v}_j^\top = \sigma^2 \sum_{j: \lambda_j > 0} \underbrace{\left(\frac{1}{(1 + \lambda/\lambda_j)^2} \right)}_{\leq 1} \cdot \frac{1}{\lambda_j} \mathbf{v}_j \mathbf{v}_j^\top.$$

17 / 29

Bias-variance trade-off

Very explicit bias-variance trade-off

$$\begin{aligned} \mathbf{w}_* - \mathbb{E}[\hat{\mathbf{w}}_{\text{ridge}}] &= \sum_{j=1}^d \frac{\lambda}{\lambda_j + \lambda} \langle \mathbf{v}_j, \mathbf{w}_* \rangle \mathbf{v}_j, \\ \text{cov}(\hat{\mathbf{w}}_{\text{ridge}}) &= (\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^\top \mathbb{E}[\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top] \mathbf{A} (\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I})^{-1}. \end{aligned}$$

λ	$\ \text{bias}\ $	variance
\uparrow	\uparrow	\downarrow
\downarrow	\downarrow	\uparrow

Using a similar analysis, can also reveal trade-off in excess expected square loss:

$$\mathbb{E}[\|\mathbf{A} \hat{\mathbf{w}}_{\text{ridge}} - \mathbf{b}\|_2^2] - \mathbb{E}[\|\mathbf{A} \mathbf{w}_* - \mathbf{b}\|_2^2].$$

18 / 29

Other interpretations

- Suppose we replace \mathbf{A} and \mathbf{b} with

$$\tilde{\mathbf{A}} := \underbrace{\begin{bmatrix} \mathbf{A} \\ \sqrt{\lambda} \mathbf{Q} \end{bmatrix}}_{(n+d) \times d \text{ matrix}}, \quad \tilde{\mathbf{b}} := \underbrace{\begin{bmatrix} \mathbf{b} \\ \mathbf{0} \end{bmatrix}}_{(n+d) \times 1 \text{ vector}},$$

where $\mathbf{Q} \in \mathbb{R}^{d \times d}$ is any orthogonal matrix.
That is, add d fictitious labeled data $(\sqrt{\lambda} \mathbf{q}_j, 0)$ for $j = 1, 2, \dots, d$, where $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_d$ is an orthonormal basis.

- Then $\tilde{\mathbf{A}}^\top \tilde{\mathbf{A}} = \mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I}$, and $\tilde{\mathbf{A}}^\top \tilde{\mathbf{b}} = \mathbf{A}^\top \mathbf{b}$.
- So $\hat{\mathbf{w}}_{\text{ols}}$ using $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{b}}$ is the same as $\hat{\mathbf{w}}_{\text{ridge}}$ using \mathbf{A} and \mathbf{b} .
- For a similar reason, $\hat{\mathbf{w}}_{\text{ridge}}$ has a certain Bayesian interpretation.

19 / 29

Ridge regression: key takeaways

- Always well-defined (for $\lambda > 0$)
- Behavior depends on eigenvectors/eigenvalues of $\mathbf{A}^\top \mathbf{A}$; the original “coordinate system” is not relevant here.
- Relative to $\hat{\mathbf{w}}_{\text{ols}}$ (when it exists): shrinks $\hat{\mathbf{w}}_{\text{ols}}$ along eigenvector directions by amount related to eigenvalue and λ .
- Regularization parameter λ is tuning-knob that controls bias-variance trade-off.
- Can be thought of as applying OLS to an augmented data set with “fake data” that ensures OLS is well-defined.

20 / 29

Sparse regression

Sparsity

Another form of regularization: only consider *sparse* w —i.e., w with only a small number ($\ll d$) of non-zero entries.

Other advantages of sparsity (especially relative to ridge):

- ▶ Sparse solutions easier to “interpret” (but caveats about interpreting weights from before still apply).
- ▶ Can be more efficient to evaluate $\langle w, x \rangle$ (both in terms of computing variable values and computing inner product).

22 / 29

Sparse regression methods

For any $T \subseteq \{1, 2, \dots, d\}$, let $\hat{w}_T :=$ OLS only using variables in T .

Subset selection

Brute-force strategy. Pick the $T \subseteq \{1, 2, \dots, d\}$ of size $|T| = k$ for which

$$\|A\hat{w}_T - b\|_2^2$$

is minimal, and return \hat{w}_T .

Gives you exactly what you want (for given value k).

Only feasible for very small k , since complexity scales with $\binom{d}{k}$.
(NP-hard optimization problem.)

Sparse regression methods

Forward stepwise regression

Greedy strategy. Starting with $T = \emptyset$, repeat until $|T| = k$:

Pick the $j \in \{1, 2, \dots, d\} \setminus T$ for which

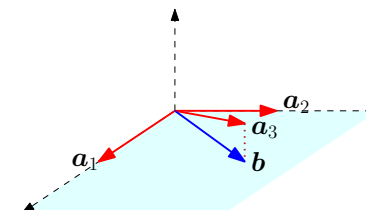
$$\|A\hat{w}_{T \cup \{j\}} - b\|_2^2$$

is minimal, and add this j to T .

Return $\hat{w}(T)$.

Gives you a k -sparse solution.

Primarily only effective when columns of A are close to orthogonal.



23 / 29

24 / 29

Aside: l_p norms

For $p \geq 1$,

$$\|v\|_p = \left(\sum_{j=1}^d |v_j|^p \right)^{1/p}.$$

In particular,

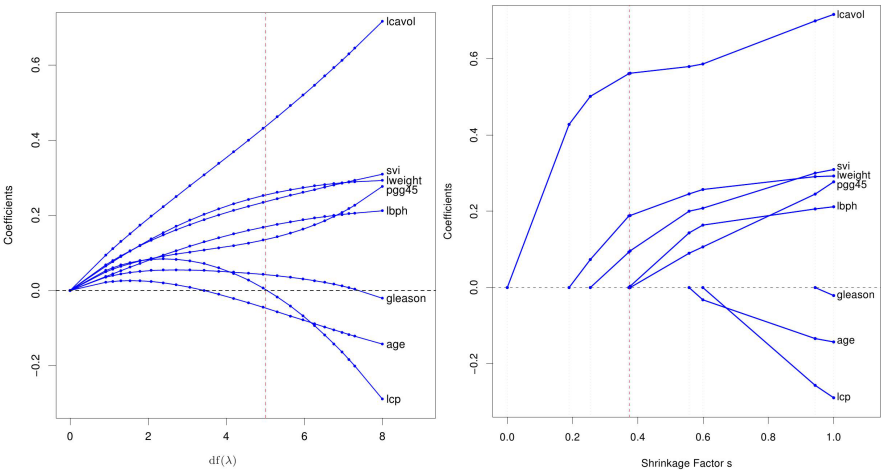
$$\|v\|_1 = \sum_{j=1}^d |v_j|.$$

These are *norms* on \mathbb{R}^d , so:

- ▶ $\|u - v\|_p$ is a valid *metric* on points in \mathbb{R}^d , and
- ▶ $\|cv\|_p = |c| \cdot \|v\|_p$ for any $v \in \mathbb{R}^d$ and $c \in \mathbb{R}$.

25 / 29

Coefficient profile (Ridge vs. Lasso)



Horizontal axis: varying λ (large λ to left, small λ to right).
Vertical axis: coefficient value in \hat{w}_{ridge} and \hat{w}_{lasso} for eight different variables.

27 / 29

Lasso (Tibshirani, 1994)

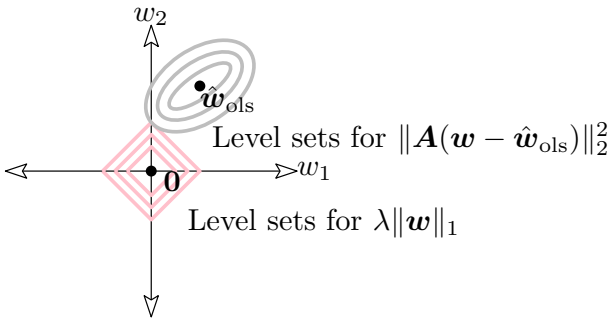
Lasso: least absolute shrinkage and selection operator

Let \hat{w}_{lasso} be a minimizer of

$$w \mapsto \arg \min_{w \in \mathbb{R}^p} \|Aw - b\|_2^2 + \lambda \|w\|_1.$$

Objective function is convex though not differentiable.

If \hat{w}_{ols} exists, then Lasso objective (as function of w) is $\|A(w - \hat{w}_{\text{ols}})\|_2^2 + \lambda \|w\|_1 + (\text{stuff not depending on } w)$.



26 / 29

Lasso: theory

Many results, mostly roughly of the following flavor.

Suppose

- ▶ $b \sim N(Aw_*, \sigma^2 I)$;
- ▶ w_* has $\leq k$ non-zero entries;
- ▶ A satisfies some special properties (typically not efficiently checkable);
- ▶ $\lambda \gtrsim \sigma \sqrt{2n \log(d)}$;

then

$$\mathbb{E}[\|w_* - \hat{w}_{\text{lasso}}\|_2^2] \leq O\left(\frac{\sigma^2 k \log(d)}{n}\right).$$

Closely related to “compressed sensing”; theory involves high-dimensional convex geometry.

28 / 29

Sparse regression: key takeaways

- ▶ **Sparsity**: a form of “regularization”, but also desirable for other reasons.
- ▶ **Subset selection**: generally intractable.
- ▶ **Greedy algorithms** (e.g., forward stepwise regression): sometimes works.
- ▶ **Lasso**: shrink coefficients towards zero in a way that tends to lead to sparse solutions.