# LINEAR REGRESSION



$R^2 = 0.06$

REXTHOR, THE DOG-BEARER
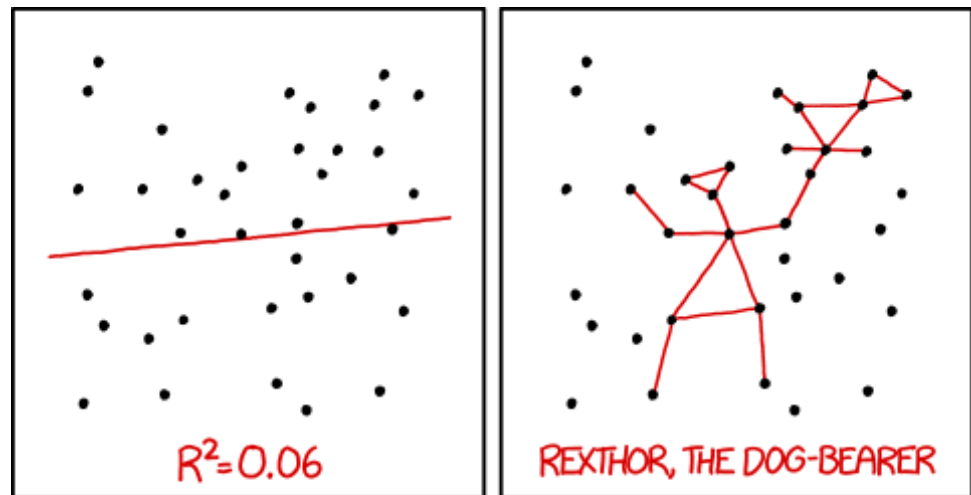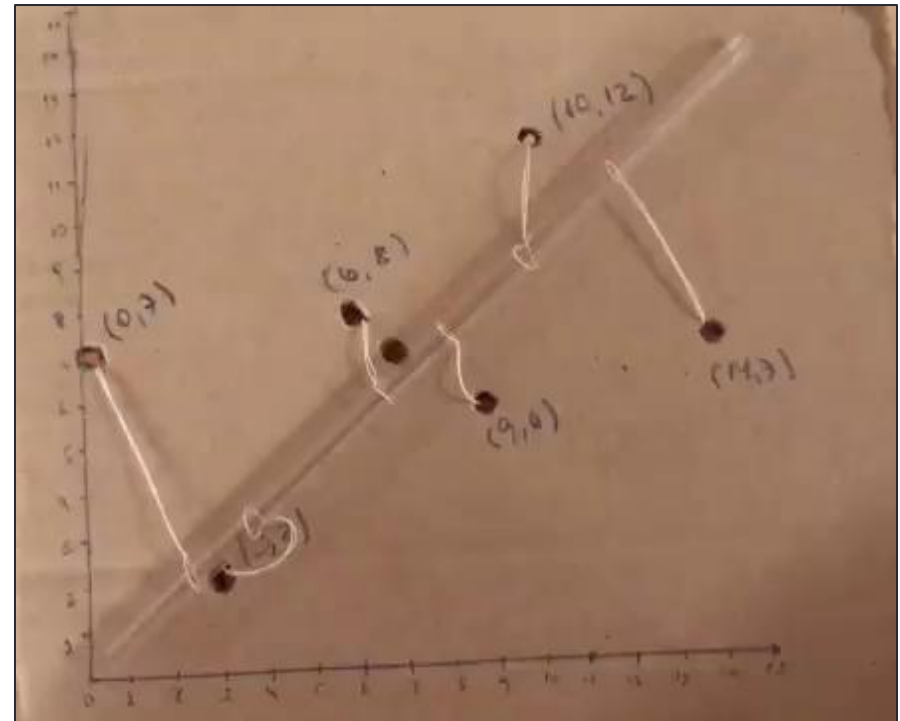
I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

Itay Evron & Ronen Nir

# Tutorial outline

- Linear regression

- The least squares problem

- Regularization

  - Ridge regression ($\ell2$ regularization)

  - LASSO ($\ell1$ regularization)

# Linear regression

- We assume $y \sim x$ and ask how

$$x \text{ explains } y$$

- Often, we assume a linear approximation $y \approx (w^*)^\top x + \epsilon$ for an unknown "ground truth" $w^* \in \mathbb{R}^d$

- Assuming a linear connection limits the search space

- We wish to find a good coefficient vector $w$
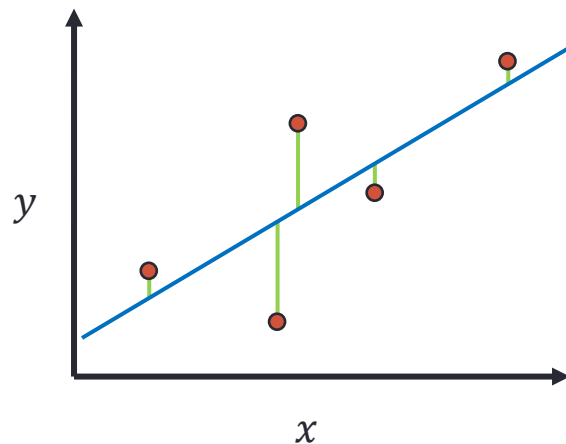


Source: @jorge_pacheco

# Optimality criterion (loss)

- For many reasons, we choose the squared error

$$\left(y_i - h(\boldsymbol{x}_i)\right)^2 = (y_i - \boldsymbol{w}^\top \boldsymbol{x}_i)^2$$

- The minimized loss:

$$\mathcal{L} = \frac{1}{m} \sum_{i=1}^{m} \left(y_i - h(\boldsymbol{x}_i)\right)^2$$

Extra – another criterion:

(Ordinary) Least Squares

Total Least Squares
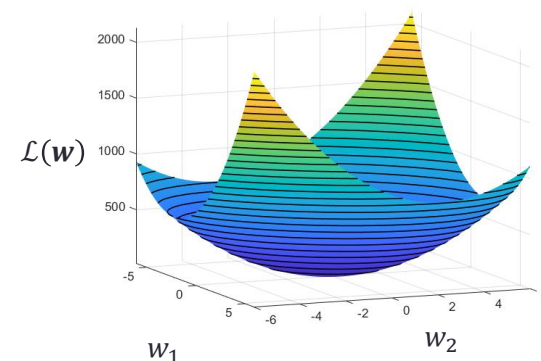


Further reading: OLS vs. TLS

# The least squares problem

- Define the squared loss over the residuals:

$$\mathcal{L}(\boldsymbol{w}) = \frac{1}{m}\sum_{i=1}^{m}(y_i - \underbrace{\boldsymbol{w}^\top \boldsymbol{x}_i}_{h(\boldsymbol{x}_i)})^2 = \frac{1}{m}\|\mathbf{X}\boldsymbol{w} - \boldsymbol{y}\|_2^2$$
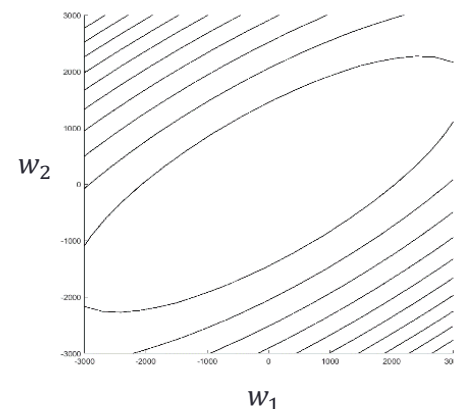
where $\mathbf{X} = \begin{bmatrix} -\boldsymbol{x}_1^\top - \\ \vdots \\ -\boldsymbol{x}_m^\top - \end{bmatrix} \in \mathbb{R}^{m\times d}$ and $\boldsymbol{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m$.

Loss landscape
of a 2d LS problem



$\mathcal{L}(\boldsymbol{w})$

$w_1$     $w_2$

- The optimization problem: $\widehat{\boldsymbol{w}} = \underset{\boldsymbol{w}}{\operatorname{argmin}} \frac{1}{m}\|\mathbf{X}\boldsymbol{w} - \boldsymbol{y}\|_2^2$

Loss level sets



$w_2$

$w_1$

- The gradient of the loss: $\nabla_{\boldsymbol{w}}\mathcal{L}(\boldsymbol{w}) = 2\mathbf{X}^\top(\mathbf{X}\boldsymbol{w} - \boldsymbol{y})$
- The Hessian of the loss: $\nabla_{\boldsymbol{w}}^2\mathcal{L}(\boldsymbol{w}) = 2\mathbf{X}^\top\mathbf{X} \succcurlyeq \mathbf{0}$
- $\Longrightarrow$ The objective loss is convex in $\boldsymbol{w}$!

# Solving least squares problems

- Derive the normal equation:

$$\nabla_{\boldsymbol{w}} \mathcal{L}(\boldsymbol{w}) = 2\mathbf{X}^\top(\mathbf{X}\boldsymbol{w} - \boldsymbol{y}) = 0 \implies \mathbf{X}^\top\mathbf{X}\widehat{\boldsymbol{w}} = \mathbf{X}^\top\boldsymbol{y}$$
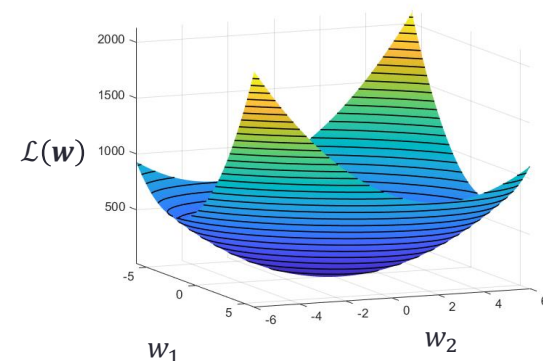
- Closed-form solution:

  - If $\mathbf{X}^\top\mathbf{X} \succ \mathbf{0}$, the unique solution is $\widehat{\boldsymbol{w}} = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\boldsymbol{y}$

  - More generally $\widehat{\boldsymbol{w}} = (\mathbf{X}^\top\mathbf{X})^+\mathbf{X}^\top\boldsymbol{y} = \mathbf{V}\boldsymbol{\Sigma}^+\mathbf{U}^\top\boldsymbol{y}$
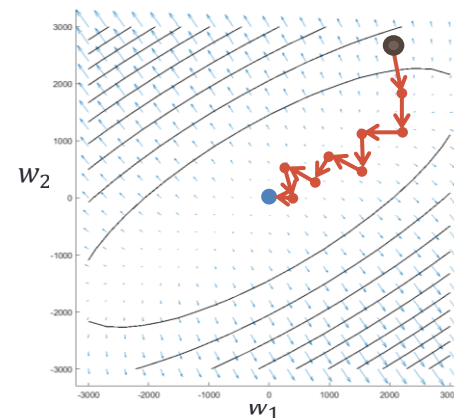
- Complexity:

  - Often, inversion is too expensive

  - Can use gradient methods $\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} - \eta\mathbf{X}^\top(\mathbf{X}\boldsymbol{w} - \boldsymbol{y})$



Loss landscape of a 2d LS problem

$\mathcal{L}(\boldsymbol{w})$

$w_1$

$w_2$



Loss level sets

$w_2$

$w_1$
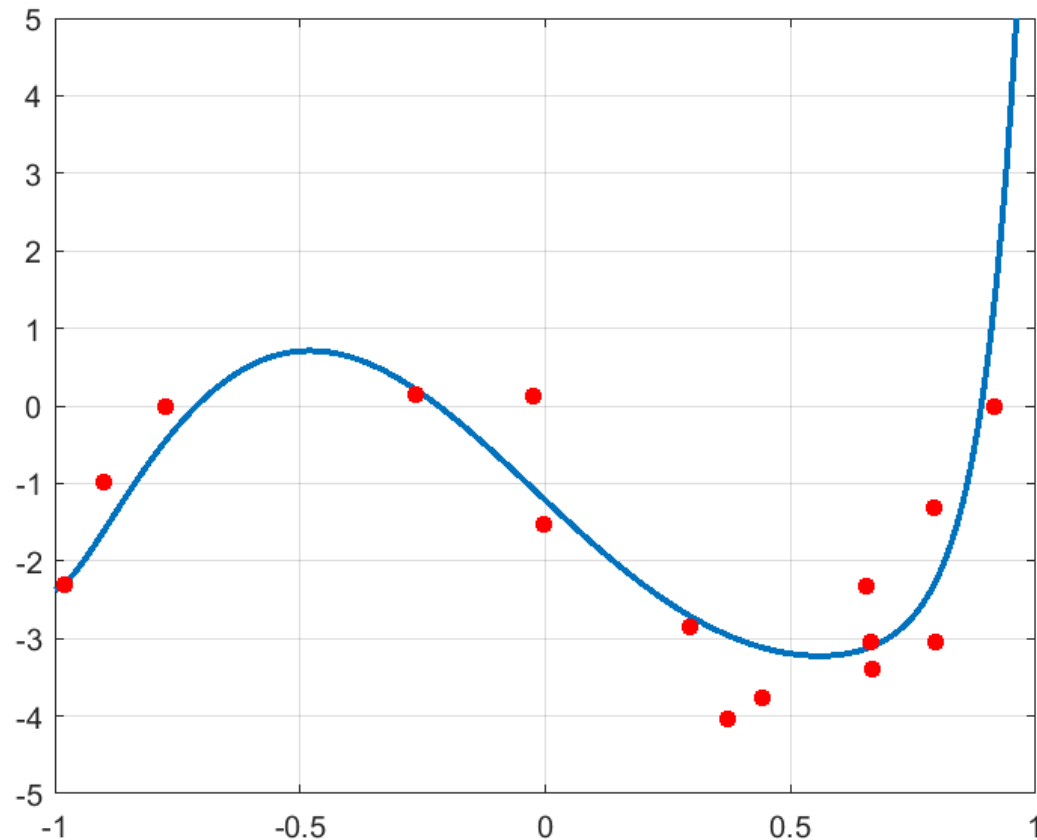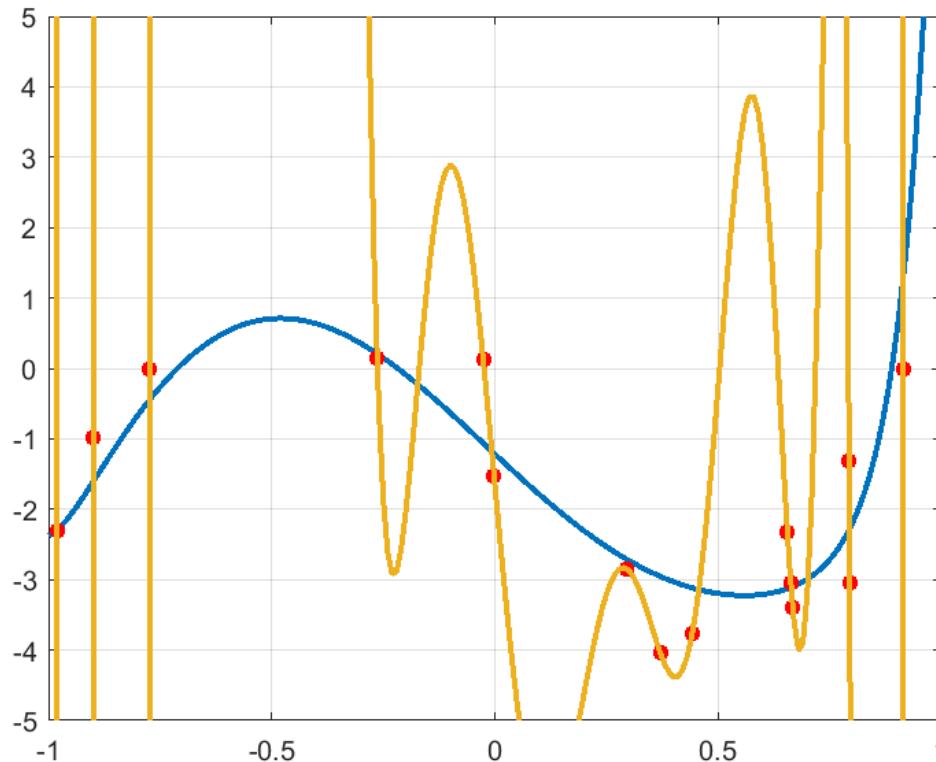
# Demo: Ridge regression for polynomial fitting

- Consider the illustrated polynomial function $f(x)$
- We sample points $(x_i, y_i)$ where $y_i = f(x_i) + \epsilon_i$ for some i.i.d noise

# Demo: Ridge regression for polynomial fitting

- We will try to fit a polynomial function of degree 25



Vandermonde matrix as a polynomial mapping:

$$X = \begin{bmatrix} - X_1^T - \\ \vdots \\ - X_m^T - \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$$

$$X' = \begin{bmatrix} x_1 & x_1^2 & \cdots & x_1^{25} \\ \vdots & & & \vdots \\ x_m & \cdots & & x_m^{25} \end{bmatrix} = \begin{bmatrix} -\varphi(x_1)- \\ \vdots \\ -\varphi(x_m)- \end{bmatrix}$$

$$\min \|X'w - y\|_2^2$$

- This is the solution that minimizes $\|\mathbf{X}w - y\|_2^2$

- Is this a good solution?

# Ridge regression ($\ell 2$ regularization)

- Regularize solutions with the $\ell 2$ norm:

$$\widehat{\boldsymbol{w}} = \operatorname*{argmin}_{\boldsymbol{w}} \left( \frac{1}{m} \sum_{i=1}^{m} (y_i - \boldsymbol{w}^\top \boldsymbol{x}_i)^2 + \lambda \|\boldsymbol{w}\|_2^2 \right) = \operatorname*{argmin}_{\boldsymbol{w}} \underbrace{\left( \frac{1}{m} \|\mathbf{X}\boldsymbol{w} - \boldsymbol{y}\|_2^2 + \lambda \|\boldsymbol{w}\|_2^2 \right)}_{\mathcal{L}_\lambda(\boldsymbol{w})}$$

- Also called Tikhonov regularization or weight decay (esp. in deep learning).

- The updated gradient and normal equation:

$$\nabla_{\boldsymbol{w}} \mathcal{L}_\lambda(\boldsymbol{w}) = \frac{2}{m} \mathbf{X}^\top (\mathbf{X}\boldsymbol{w} - \boldsymbol{y}) + 2\lambda \boldsymbol{w} \implies \underbrace{(\mathbf{X}^\top \mathbf{X} + m\lambda \mathbf{I}_{d \times d})}_{>0} \widehat{w} = \mathbf{X}^\top \boldsymbol{y}$$
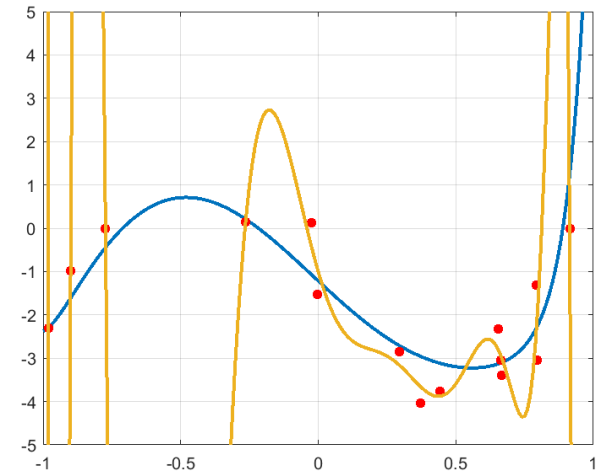
- The shrinkage effect
  - We should understand what happens in the limits $\lambda \to 0$ and $\lambda \to \infty$

- Optimization
  - We can compute the closed form solution $\widehat{\boldsymbol{w}} = (\mathbf{X}^\top \mathbf{X} + m\lambda \mathbf{I}_{d \times d})^{-1} \mathbf{X}^\top \boldsymbol{y}$
  - Loss remains convex and differentiable, so we can still run GD

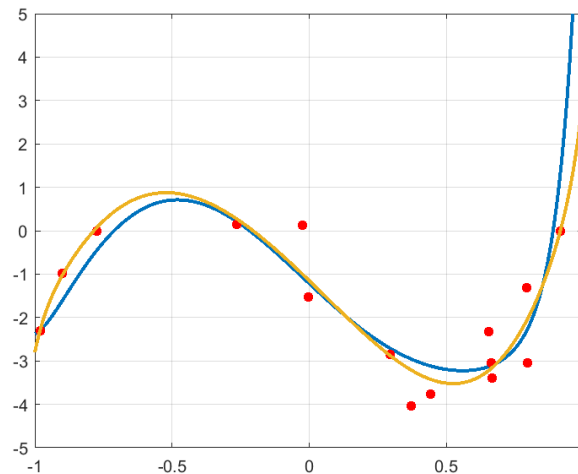# Demo: Ridge regression for polynomial fitting
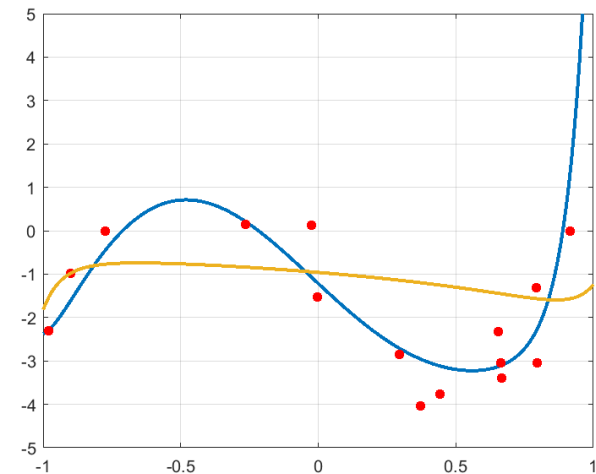
$\lambda = 0$

$\lambda = 10^{-8}$

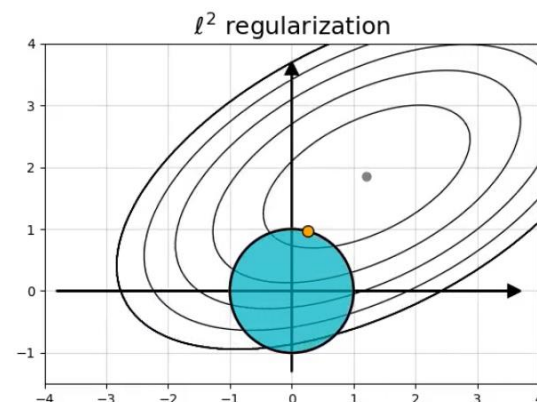Regularization mitigates overfitting and helps generalization!

$\lambda = 10^{-2}$

$\lambda = 10$

# Equivalence to constrained problems

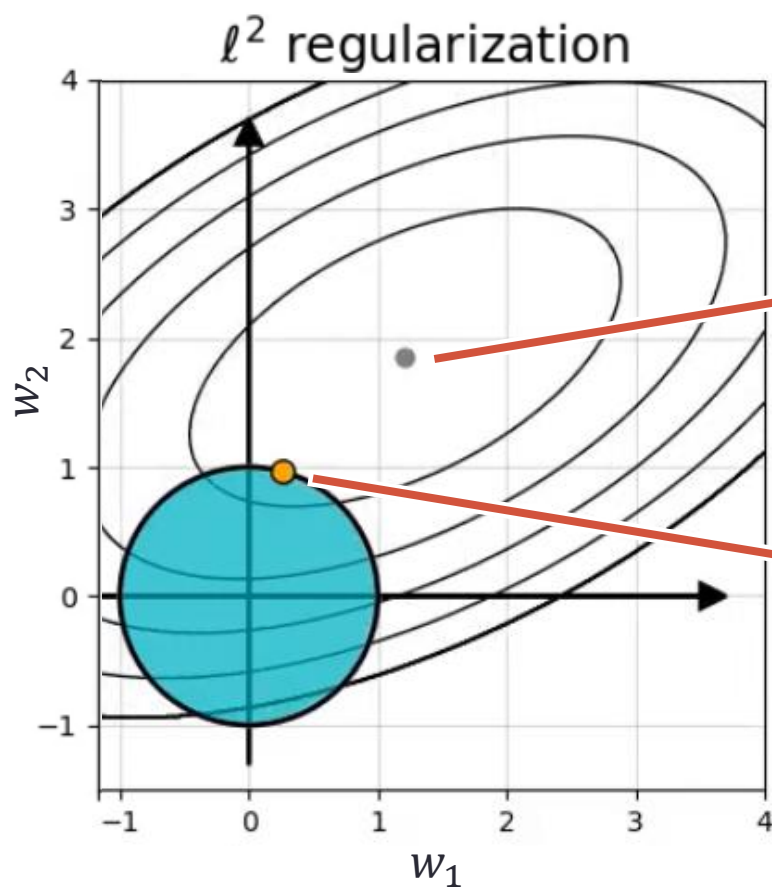- Theorem: the regularized problems are equivalent to unregularized problems with norm constraints.

$$\boldsymbol{w}^{\text{Ridge}} = \underset{\boldsymbol{w}}{\operatorname{argmin}}\left(\frac{1}{m}\Sigma_{i=1}^{m}(y_i - \boldsymbol{w}^{\top}\boldsymbol{x}_i)^2 + \lambda\|\boldsymbol{w}\|_2^2\right)$$

$$= \underset{\boldsymbol{w}}{\operatorname{argmin}}\frac{1}{m}\Sigma_{i=1}^{m}(y_i - \boldsymbol{w}^{\top}\boldsymbol{x}_i)^2, s.t. \ \|\boldsymbol{w}\|_2^2 \le c$$
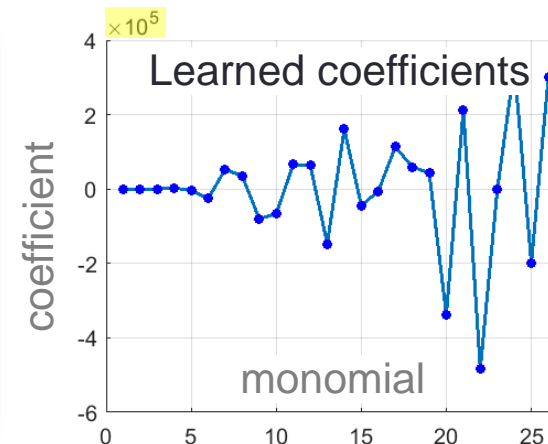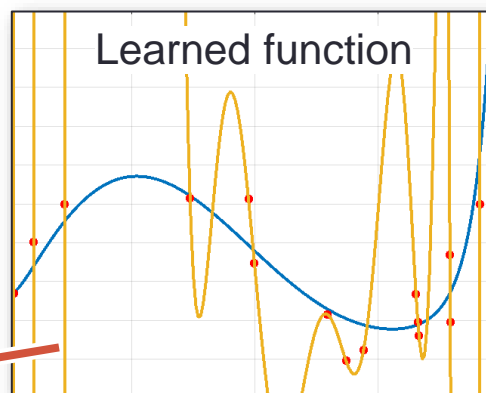


$\ell^2$ regularization

where $\lambda, c$ are related.

# Understanding the solution space



Minimal train loss solution (unregularized solution)

Bounded norm solution (regularized solution)

(illustration in 2d)

# LASSO ($\ell 1$ regularization)

- Least Absolute Shrinkage and Selection Operator

- Regularize solutions with the $\ell 1$ norm:

$$\widehat{\boldsymbol{w}} = \underset{\boldsymbol{w}}{\text{argmin}} \left( \frac{1}{m} \sum_{i=1}^{m} (y_i - \boldsymbol{w}^\top \boldsymbol{x}_i)^2 + \lambda \|\boldsymbol{w}\|_1 \right)$$

- Often induces sparse solutions (few nonzero entries)

- No closed-form solution!

- Optimization

  - Loss remains convex but no longer differentiable.

  - Could run subgradient descent, but more suitable algorithms exist.

    - Think: how will the subgradient method perform on $f(x) = |x|$ ?

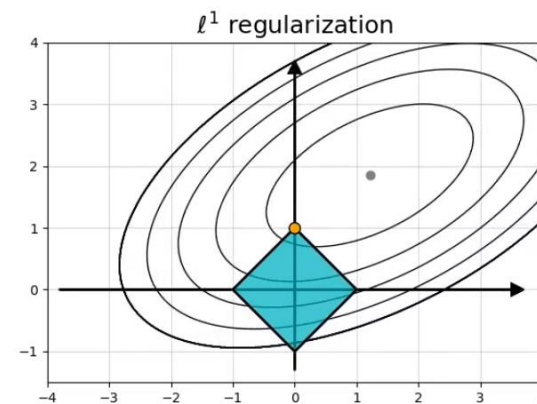  - Further reading: Why not vanilla GD?, Proximal gradient methods for learning, FISTA

# Equivalence to constrained problems

- Theorem: the regularized problems are equivalent to unregularized problems with norm constraints.

$$w^{\text{Ridge}} = \underset{w}{\text{argmin}}\left(\frac{1}{m}\Sigma_{i=1}^{m}(y_i - w^\top x_i)^2 + \lambda\|w\|_2^2\right)$$

$$= \underset{w}{\text{argmin}}\frac{1}{m}\Sigma_{i=1}^{m}(y_i - w^\top x_i)^2 , s.t. \ \|w\|_2^2 \leq c$$



$\ell^2$ regularization

$$w^{\text{LASSO}} = \underset{w}{\text{argmin}}\left(\frac{1}{m}\Sigma_{i=1}^{m}(y_i - w^\top x_i)^2 + \lambda\|w\|_1\right)$$

$$= \underset{w}{\text{argmin}}\frac{1}{m}\Sigma_{i=1}^{m}(y_i - w^\top x_i)^2 , s.t. \ \|w\|_1 \leq c$$
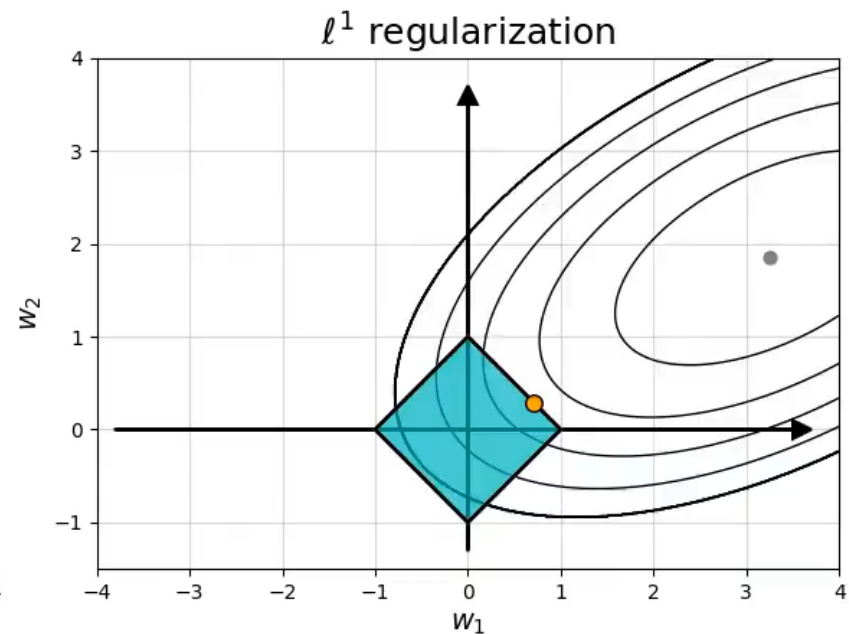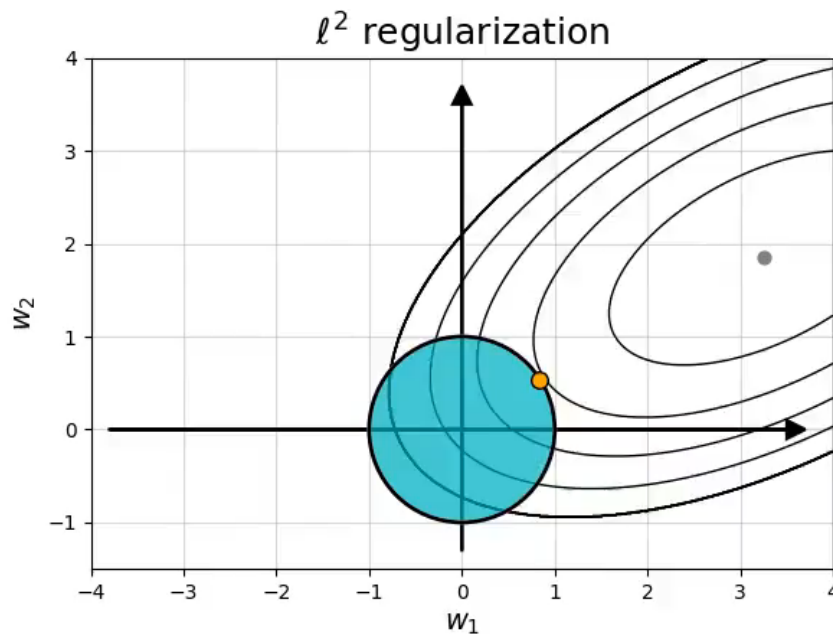
where $\lambda, c$ are related.



$\ell^1$ regularization

- Illustrate why LASSO induces sparser models.

# LASSO induces sparse models



$\ell^1$ induces sparse solutions for least squares

by @itayevron

The level sets belong to an unregularized least squares problem.
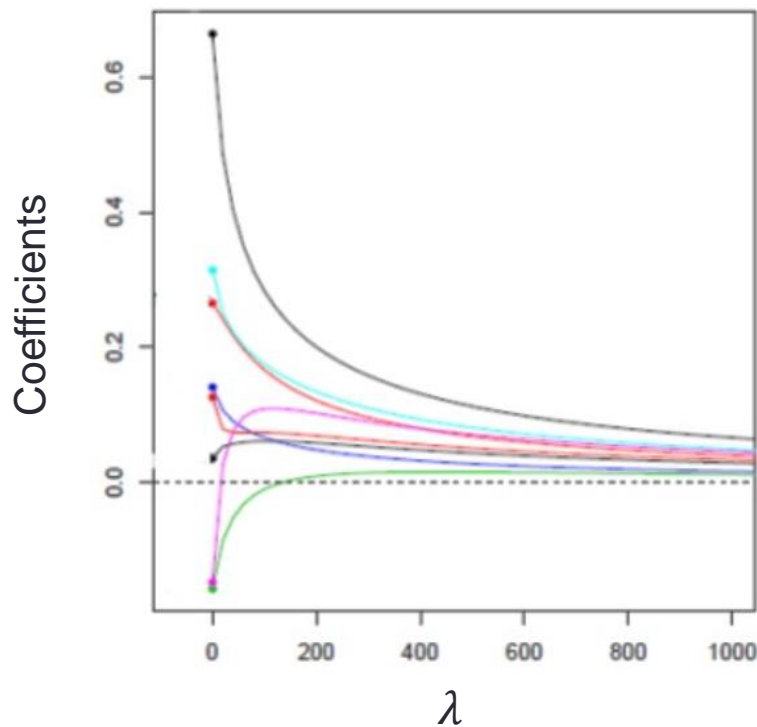The orange points have the lowest LS error on each unit "circle".
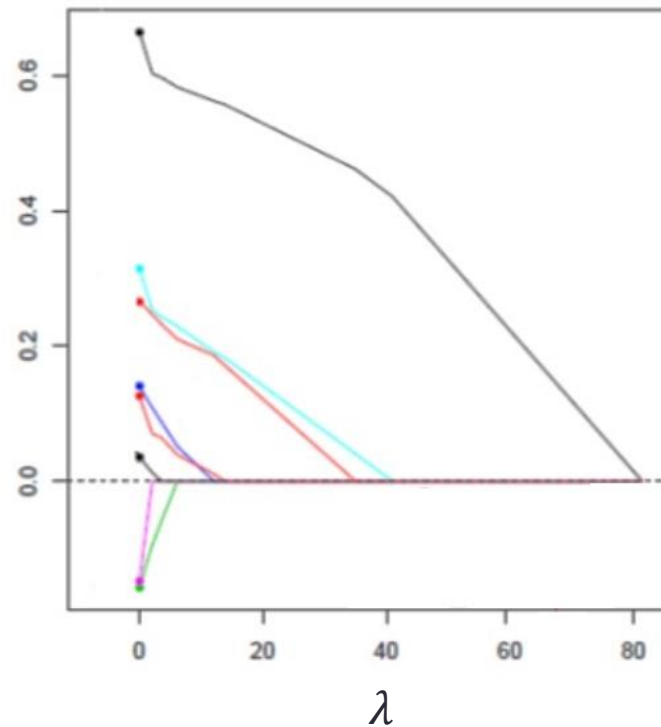Animation can be found on GitHub.

Notice: in both cases we don't get the solution with the minimal (unregularized) training MSE!

# Different regularization behaviors

$\ell^2$ regularization
causes weight decay

$\ell^1$ regularization causes
"variable selection"



Extra:  Are larger <u>unregularized</u> coefficients necessarily more "important"?

Answer:  Not necessarily! See Q3 in Exam A of Winter 2020-21

# Summary

- Linear regression tries to linearly "explain" labels $y$ using feature vectors $x$.

- Often formulated by least squares.

- Regularization can help prevent overfitting.

- Different regularizations induce different solutions.