

CS 189/289

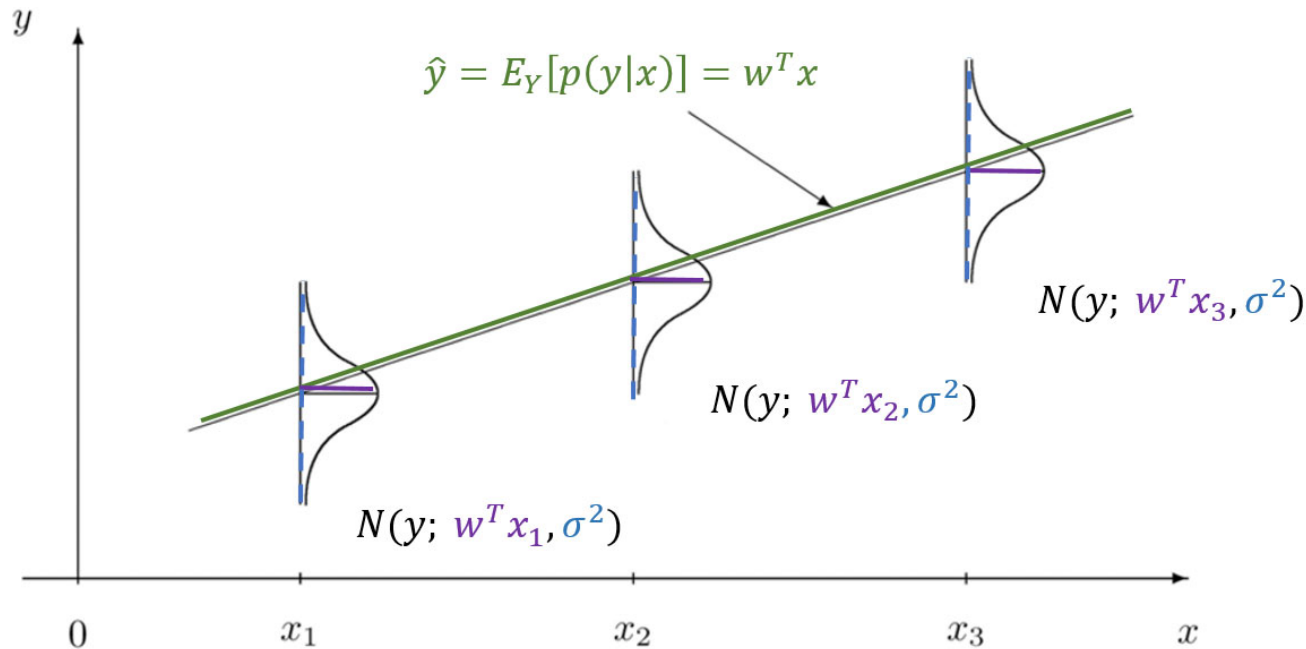
Today's lecture:

Linear regression part II

Reloading last lecture

$$\begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} = A$$

Gaussian linear regression, $p(y|x) = N(y|w^T x, \sigma^2)$



- $\theta_{MLE} = (w_{MLE}, \sigma_{MLE}^2) = \arg \max_{(w, \sigma^2)} \log p(D|\theta)$
- $\mathcal{L}_w = (y - Aw)^T (y - Aw)$, set $\frac{\partial \mathcal{L}}{\partial w} = 0$
- $A^T y = A^T A w$ ($A \in \mathbb{R}^{N \times d}$)
- $w_{MLE} = (A^T A)^{-1} A^T y$, $\sigma_{MLE}^2 = \frac{1}{N} \sum_i (y_i - w^T x)^2$

careful!

- When not invertible, there are ∞ many equally good solutions for w_{MLE} .
- Called *underdetermined* linear regression.

- **Not invertible** if columns (features) in A are linearly dependent.
- Automatically happens when $d > N$, in which case, $y = Aw$ exactly.

Intuition of why ∞ # of w_{MLE} solutions

- Suppose we have 2 *linearly dependent features* in the training data such that $\alpha x_1 = x_2$.
- Suppose we found one MLE solution, \hat{w} .
- Then for any training data point, $\hat{y} = x^T \hat{w}$.

$$= [x_1 \ x_2] \begin{bmatrix} \hat{w}_1 \\ \hat{w}_2 \end{bmatrix}$$

$$= \hat{w}_1 x_1 + \hat{w}_2 x_2 = \hat{w}_1 x_1 + \hat{w}_2 \alpha x_1$$

$$= (\hat{w}_1 + \hat{w}_2 \alpha) x_1 = (\hat{w}_1 + \hat{w}_2 \alpha + \beta - \beta) x_1 \text{ for any } \beta.$$

$$= ((\hat{w}_1 + \beta) + (\hat{w}_2 \alpha - \beta)) x_1 \text{ for any } \beta.$$

$$= (\hat{w}_1 + \beta) x_1 + (\hat{w}_2 \alpha - \beta) \frac{x_2}{\alpha} = (\hat{w}_1 + \beta) x_1 + (\hat{w}_2 \alpha - \beta) \frac{1}{\alpha} x_2$$

$$= [x_1 \ x_2] \begin{bmatrix} \hat{w}_1 + \beta \\ (\hat{w}_2 \alpha - \beta) / \alpha \end{bmatrix} = x^T \tilde{w}$$

Of all the $\{w_{MLE}\}$ with zero error, is there one that intuitively might be generally be better?

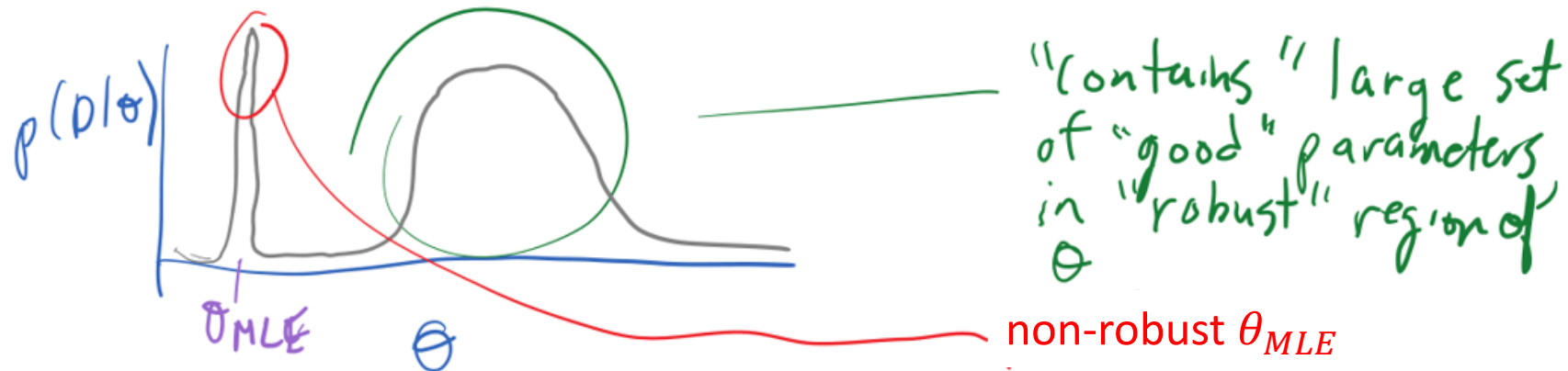
Intuition for choosing one specific $\hat{\mathbf{w}}$.

- Of the ∞ solutions for \mathbf{w}_{MLE} , choose the one with the least norm, $\|\mathbf{w}_{MLE}\|_2$. Why might this be a good idea?
- Hint 1: smaller norm tends to have smaller individual values.
- Hint 2: don't expect co-linearity of features with test data.
- Consider prediction, $\hat{\mathbf{y}} = \mathbf{w}^T \mathbf{x}$. How much does the prediction change when we perturb, $\mathbf{x}' = \mathbf{x} + \delta$, for different norm \mathbf{w} ?
- With smaller coefficients the model is less sensitive to noise.
- What about in non-degenerate linear regression ($\mathbf{A}^T \mathbf{A}$ is invertible)?
- Yes! For many problems (and models), small param norm is a good idea.
- This is one e.g. of *regularization*: in effect, reduce # free parameters, while keeping the same set of parameters!

Recall how MLE can go wrong?

MLE yields a "point estimate" of our parameter

- When we perform MLE, we get just one single estimate of the parameter, θ , rather than a distribution over it which captures uncertainty.
- In Bayesian statistics, we obtain a (posterior) distribution over θ . We will touch more on this in a few lectures.



L2 regularized linear regression

To shrink \mathbf{w} to be smaller than the MLE solution, we add a “penalty” term to the loss function:

$$\mathcal{L} = (\mathbf{y} - A\mathbf{w})^T (\mathbf{y} - A\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$
$$\mathbf{w}_{L_2} = \operatorname{argmin}_{\mathbf{w}} (\mathbf{y} - A\mathbf{w})^T (\mathbf{y} - A\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

Also called “Ridge” regression, or L2 linear regression.

Related to *Bayesian* modeling (next).

The Bayesian modelling approach



- Bayesians put a *prior distribution* on the parameters, $p(\theta)$.
- Then they seek to compute the *posterior distribution*, $p(\theta|D)$.
- Then, predictive distribution is given by

$$p(y|x) = \int_{\theta} p(y|x, \theta)p(\theta|D)d\theta = E_{\theta}[p(y|x, \theta)]$$

- Procedurally, this is done using Bayes' rule:

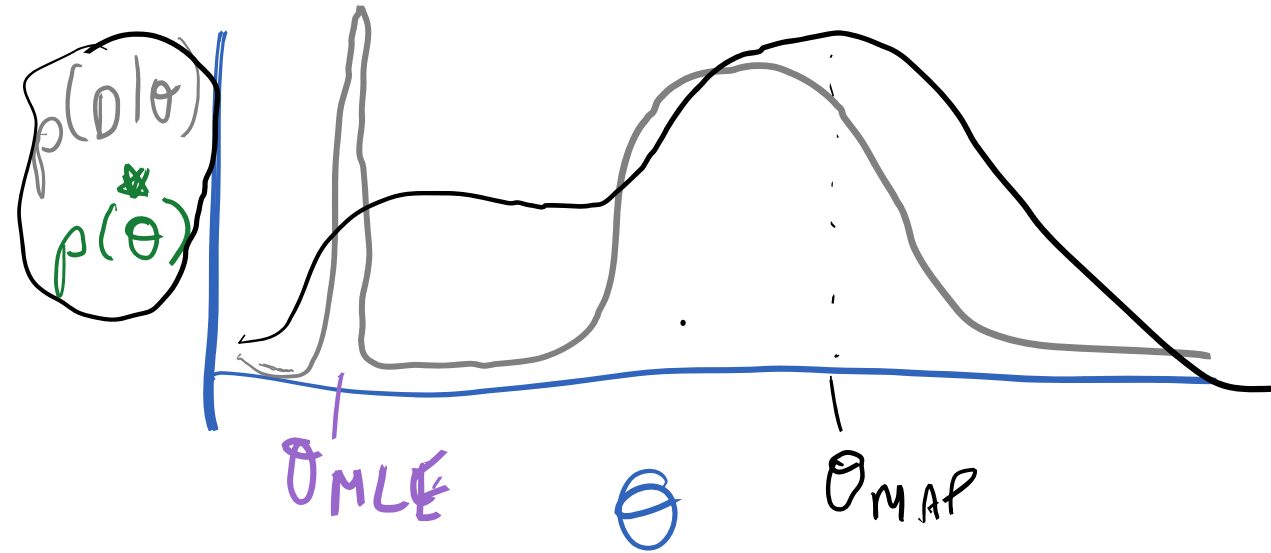
$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}.$$

- Difficult in practice! $p(D) = \int_{\theta} p(D, \theta)d\theta = \int_{\theta} p(D|\theta)p(\theta)d\theta$
- We will be lazy, instead being pseudo Bayesians, yielding L2 regression:
- $\theta_{lazy} = \operatorname{argmax}_{\theta} p(\theta|D)$ *Maximum A Posteriori* (MAP) estimation.

MAP: the lazy Bayesian (*Maximum A Posteriori*)

- Still use a prior over parameters, $p(\theta)$.
- Finds point estimate of the parameter that maximizes the posterior.
- $\theta_{MAP} = \operatorname{argmax}_{\theta} p(\theta|D)$

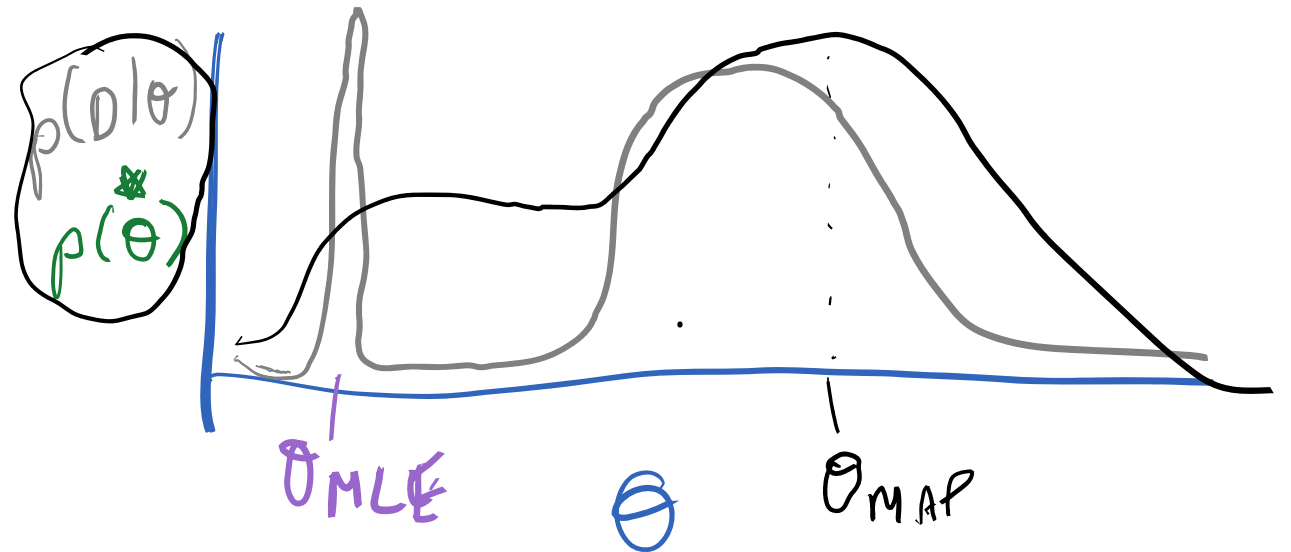
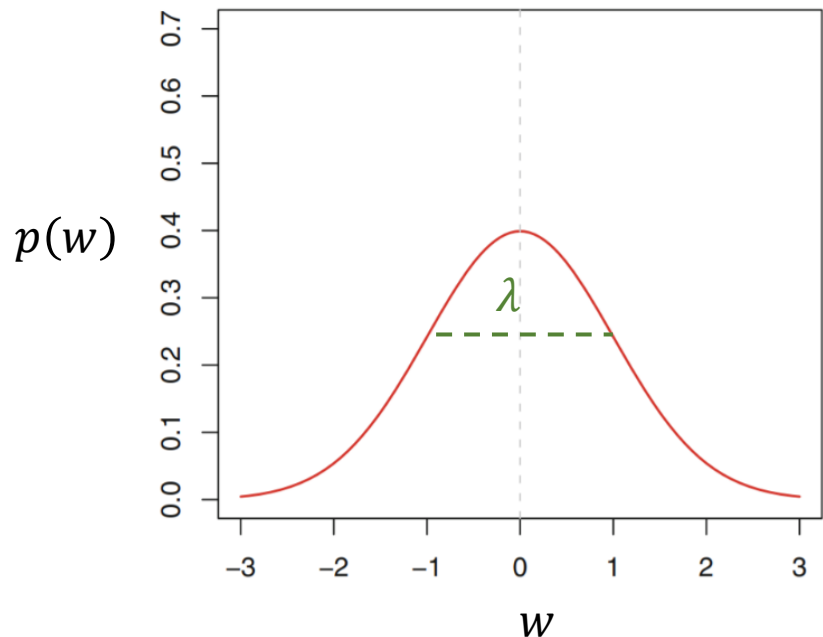
$$\begin{aligned} &= \operatorname{argmax}_{\theta} \frac{p(D|\theta)p(\theta)}{p(D)} \\ &= \operatorname{argmax}_{\theta} p(D|\theta)p(\theta) \end{aligned}$$



$$p(D) = \int_{\theta} p(D, \theta) d\theta = \int_{\theta} p(D|\theta)p(\theta) d\theta$$

A prior for small weights yields L2 regression!

- Zero-mean prior, $p(w) = N(w; 0, \lambda I)$.
- Bayesian posterior, $p(w|D) = \frac{p(D|w)p(w)}{p(D)}$ is then “nice” in that everything is Gaussian (can work it out using MVGs).



MAP for linear regression w Gaussian prior

$$\begin{aligned}w_{MAP} &= \operatorname{argmax}_w \log p(D|w) p(w) = \operatorname{argmax}_w \log p(D|w) + \log N(w; 0, \lambda I) \\&= \operatorname{argmax}_w \sum_{i=1}^N \log N(y_i | w^T x_i, \sigma^2) + \log N(w | 0, \lambda I) \\&= \operatorname{argmin}_w \frac{1}{2\sigma^2} (y - Aw)^T (y - Aw) - \sum_{i=1}^d \log \left[\frac{1}{\sqrt{2\pi\lambda}} \exp \left(-\frac{(w_i - 0)^2}{2\lambda} \right) \right] \\&= \operatorname{argmin}_w \frac{1}{2\sigma^2} (y - Aw)^T (y - Aw) + \sum_{i=1}^d \left[-\log \frac{1}{\sqrt{2\pi\lambda}} + \frac{w_i^2}{2\lambda} \right] \\&= \operatorname{argmin}_w \frac{1}{2\sigma^2} (y - Aw)^T (y - Aw) + \sum_d \frac{1}{2\lambda} w_i^2 \\&= \operatorname{argmin}_w \frac{1}{2\sigma^2} (y - Aw)^T (y - Aw) + \frac{1}{2\lambda} \|w\|_2^2 \\&= \operatorname{argmin}_w (y - Aw)^T (y - Aw) + 2\sigma^2 \frac{1}{2\lambda} \|w\|_2^2 \\&= \operatorname{argmin}_w (y - Aw)^T (y - Aw) + \lambda' \|w\|_2^2 \quad \text{for } \lambda' = \frac{\sigma^2}{\lambda}.\end{aligned}$$

Equivalence between
MAP w Gaussian prior
and L2 regression!

Obtaining the MAP/L2 solution

$$\begin{aligned}w_{L_2} &= \operatorname{argmin}_w (y - Aw)^T (y - Aw) + \lambda \|w\|_2^2 \\ &= \operatorname{argmin}_w (y - Aw)^T (y - Aw) + \lambda w^T w\end{aligned}$$

Take partial derivative and set to zero:

$$\nabla_w \mathcal{L}_{MAP} = -2A^T y + 2A^T A w + 2\lambda I w$$

$$\rightarrow 0 = -A^T y + A^T A w + \lambda I w$$

$$\rightarrow A^T y = (A^T A + \lambda I) w$$

$$\rightarrow (A^T A + \lambda I)^{-1} A^T y = w$$

$$\text{So } w_{L_2} = (A^T A + \lambda I)^{-1} A^T y.$$

If $\lambda > 0$, we can invert $(A^T A + \lambda I)$.

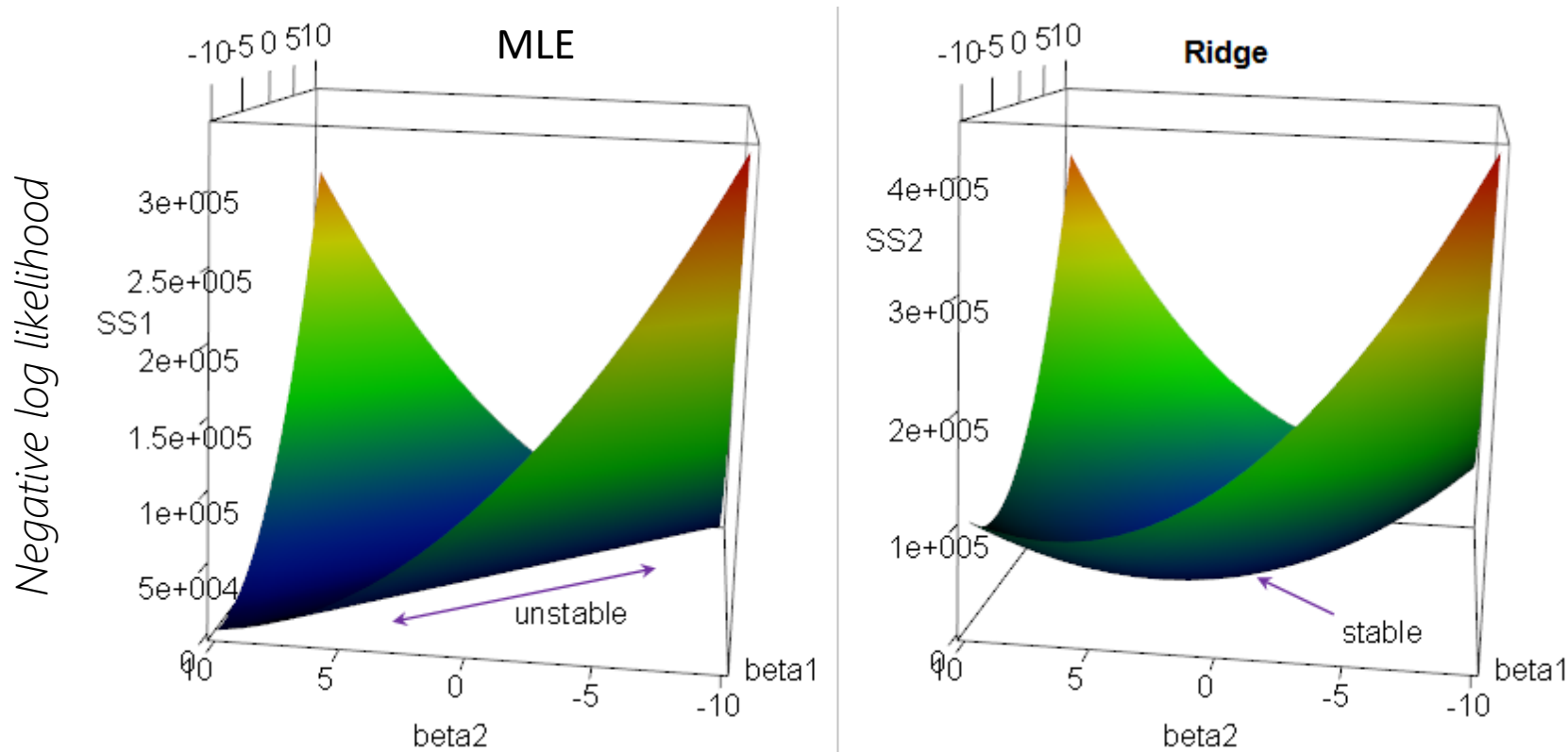
$$A = \Phi D \Phi^T$$

$$A^{-1} = \Phi D^{-1} \Phi^T$$

$$\text{where } D^{-1} = \begin{bmatrix} \lambda_1^{-1} & & & \\ & \lambda_2^{-1} & & \\ & & \ddots & \\ & & & \lambda_n^{-1} \end{bmatrix}$$

Aside, why is this called "Ridge" regression?

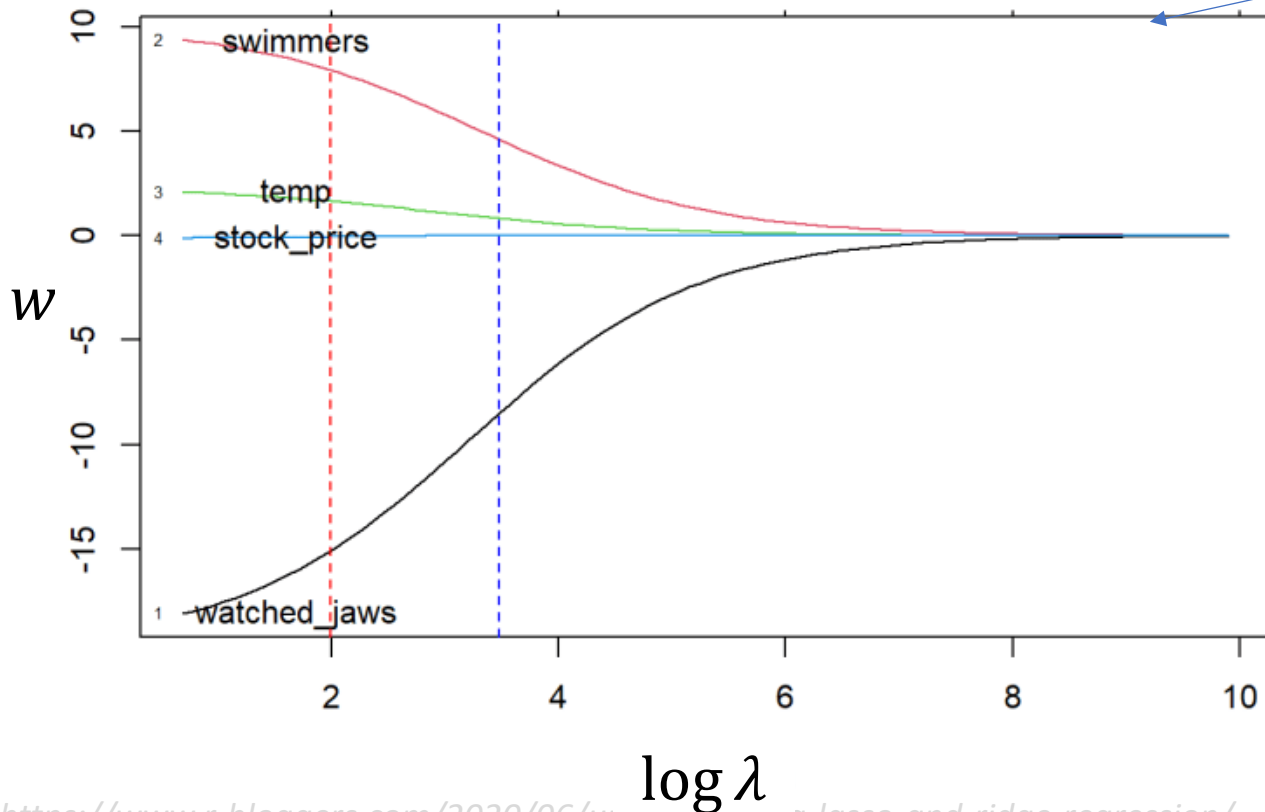
When some features are linearly dependent (can't invert $A^T A$), we have ∞ many equally good solutions that form a ridge.



Effect of value of λ

$$\mathcal{L}_{MAP} = (y - Aw)^T (y - Aw) + \lambda \|w\|_2^2$$

of non-zero features

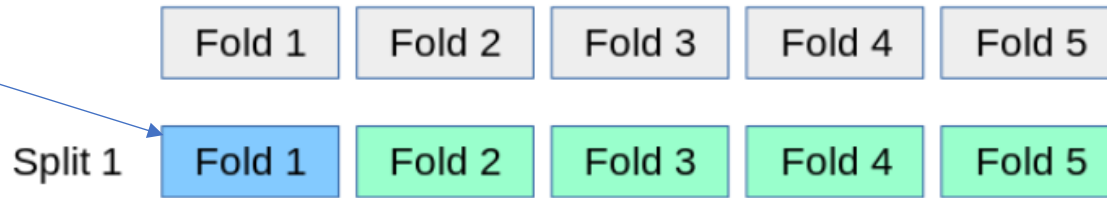


- Practically, how should we set λ ?
- Can we treat it as a parameter in the loss, and minimize wrt it?
- No: cannot use MLE!
- Need independent data, a *validation set* on which to evaluate the loss.

Train/validation/test split



Validation set from
the training data



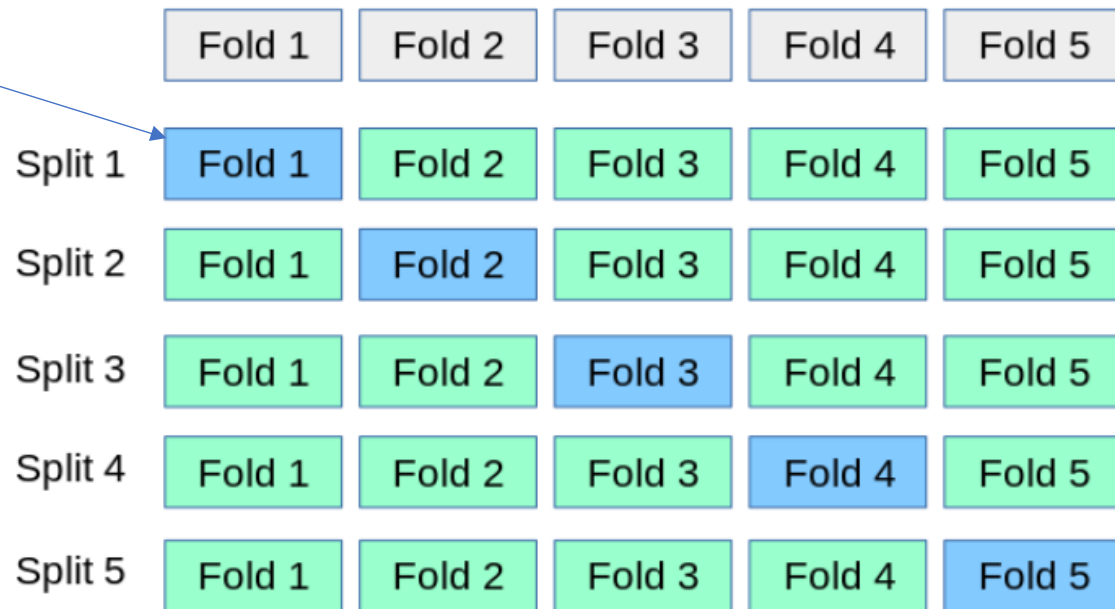
1. Find value of hyperparam that is best on the **validation set**.
2. Asses performance on the **test data**.

How to assess? Compute the log likelihood of the validation/train data (so also estimate $\hat{\sigma}^2$).

K-fold cross-validation



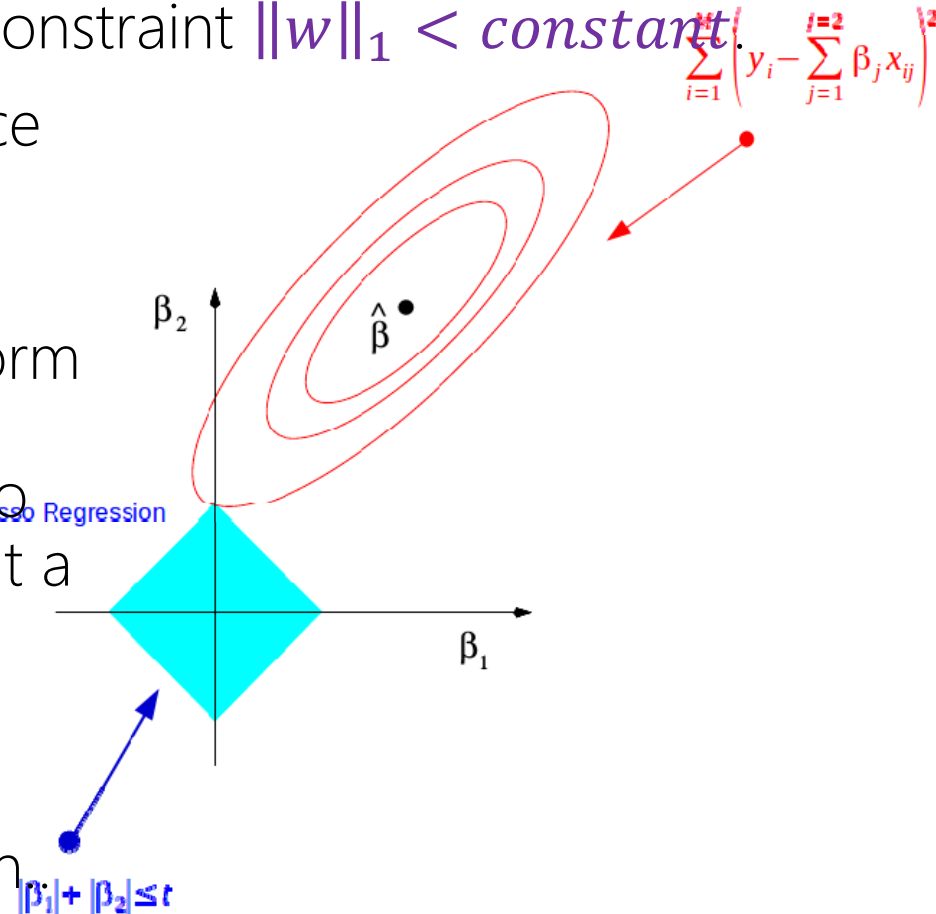
Validation set from the training data



1. Find value of hyperparam that is best across all **validation sets**.
2. Asses performance on the test data.

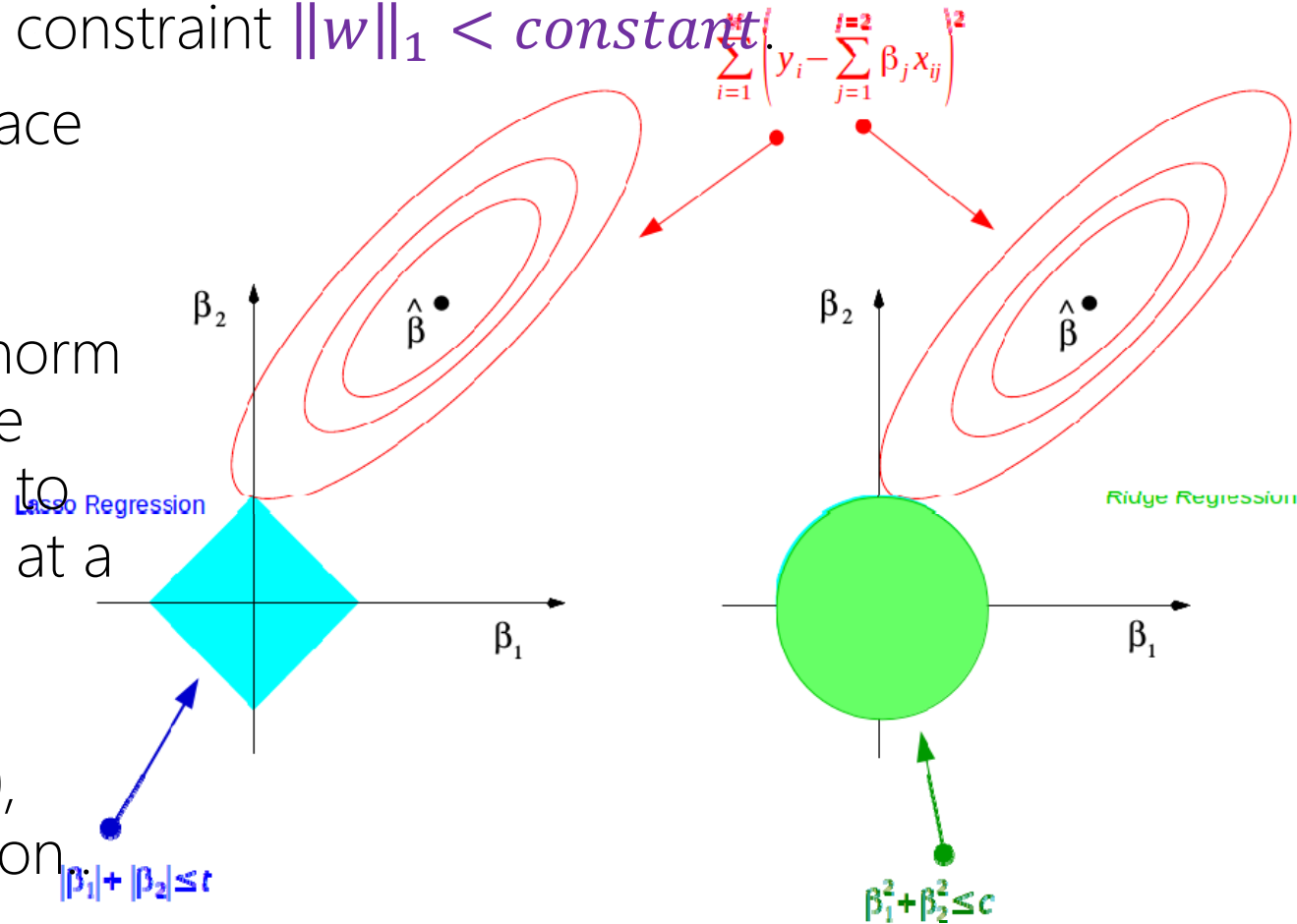
L_1 -penalized linear regression, aka *Lasso*

- $w_{L_1} = \underset{w}{\operatorname{argmin}} (y - Aw)^T (y - Aw) + \lambda \|w\|_1$
- Why does the L_1 norm penalty tends to induce sparse w ?
- Equivalent to MLE with constraint $\|w\|_1 < \text{constant}$.
- "Pointy" constraint surface is jutting out along the axes.
- In many cases, the L_1 norm constraint will cause the unconstrained solution to intersect the constraint at a corner.
- The corners are where some coefficients are 0, which is a sparse solution.

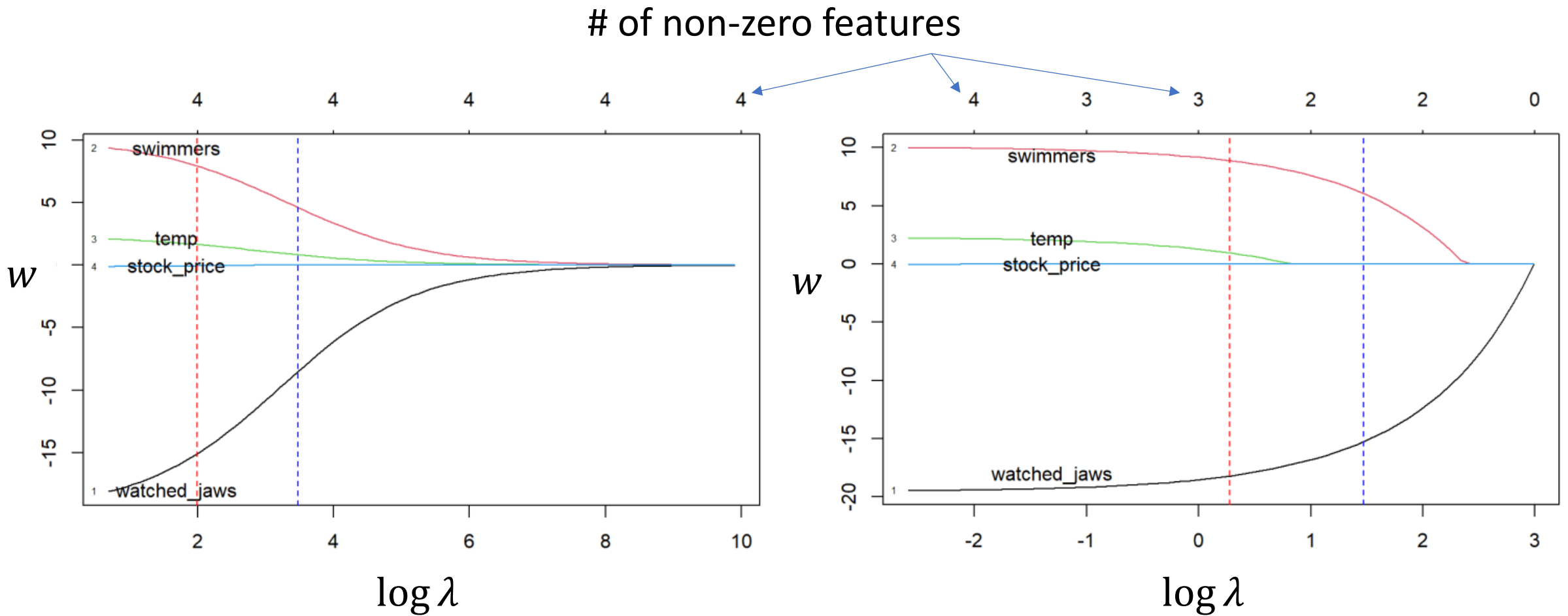


L_1 -penalized linear regression, aka *Lasso*

- $w_{L_1} = \operatorname{argmin}_w (y - Aw)^T (y - Aw) + \lambda \|w\|_1$
- Why does the L_1 norm penalty tends to induce sparse w ?
- Equivalent to MLE with constraint $\|w\|_1 < \text{constant}$.
- “Pointy” constraint surface is jutting out along the axes.
- In many cases, the L_1 norm constraint will cause the unconstrained solution to intersect the constraint at a corner.
- The corners are where some coefficients are 0, which is a sparse solution.

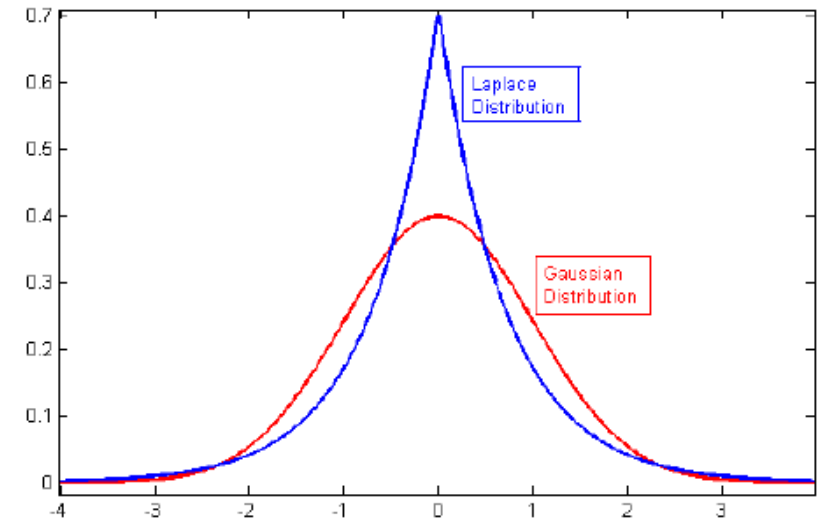


Ridge vs Lasso: shrinkage vs sparsity



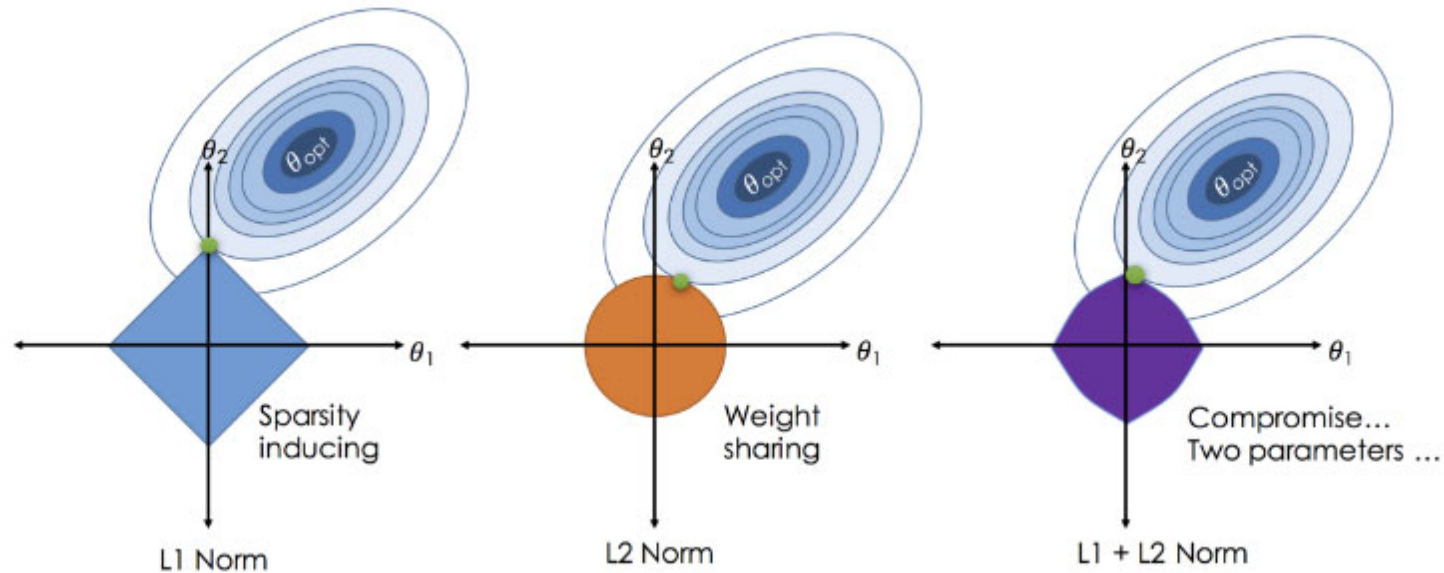
MAP interpretation for Lasso/ L_1 -penalized linear regression?

- L_2 regression arose from a $N(\mathbf{0}, \lambda I)$ prior.
- Is there a prior corresponding to L_1 ?
- Technically, the Laplace prior, $p(\mathbf{w}) = \exp(-\lambda' \|\mathbf{w}\|_1)$.



Combine L_1 and L_2 penalties?

Yes, "elastic net regression".



Issues with LASSO:

- When $d \gg N$, will select no more than N features.
- If highly correlated features, tends to ignore all but one.