

Today (Linear Regression Ctd)

- Statistical intuition on Ridge and LASSO
(distinction between MLE (Maximum Likelihood) and the MAP (Maximum A Posteriori estimator))
- Comparison between Ridge and LASSO through constrained formulations (l_1 balls)
- Effect of regularization (Ridge & LASSO) on bias and variance
- Plot bias and variance decomposition for the MSE

→ Classification

→ Simple least squares binary classifier

↳ Extensions to Multiple classes setting

→ logistic regression

Statistical Motivation for Ridge and LASSO

$$l_{\text{Ridge}}(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - \beta^T \tilde{X}^{(i)})^2 + \lambda \sum_{j=1}^D |\beta_j|^2$$

$$\tilde{X}^{(i)} \in \mathbb{R}^{D+1}$$

$$\tilde{X}^{(i)} = [1, X^{(i)}]$$

$$\{t^{(i)}, X^{(i)}\}_{i=1}^N$$

training data

training data

Assumption 1

$$t^{(i)} = \beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)} + \varepsilon^{(i)}$$

Assumption 2

$$\varepsilon^{(i)} \sim \mathcal{N}(0, \sigma)$$

$\mathcal{N}(0, \sigma) = \text{Gaussian distr.}$

$$\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Assumption 3

$\varepsilon^{(i)}$ independent

(i.i.d independent & identically distributed)

Since the only random contribution to $t^{(i)}$ is $\varepsilon^{(i)}$ the probability of observing a particular $t^{(i)}$ can be deduced from $p(\varepsilon^{(i)})$. In particular

$$\varepsilon^{(i)} = t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})$$

$$p(t^{(i)} | x^{(i)}; \beta) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ - \frac{(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2}{2\sigma^2} \right\}$$

→ Now assume that we have some prior intuition on how the β_j (parameters we want to estimate) should look like.

and assume that this intuition is given to us in the form of a probability distribution $p(\beta_j)$

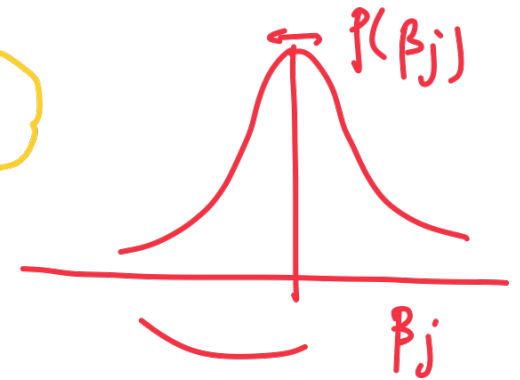
Further assume that β_j 's are independent.

(Gaussian with 0 variance λ^2)

As an example let us take

$$p(\beta_j) = \frac{1}{\sqrt{2\pi}\lambda} \exp\left(-\frac{(\beta_j)^2}{2\lambda^2}\right) \quad (*)$$

$$p(\beta) = \prod_{j=1}^n p(\beta_j)$$



Question? How can we incorporate this

additional prior information on β_j 's in

order to improve our MLE estimator?

Posterior distribut (*)

→ Use Bayes' theorem

$$p(\beta | t^{(i)}) = \frac{p(t^{(i)} | \beta) p(\beta)}{p(t^{(i)})} \quad (*)$$

$$p(t^{(i)}) = \sum_{\beta_k} p(t^{(i)} | \beta_k) p(\beta_k) \rightarrow \text{constant across different choices of } \beta_k$$

A good estimator for β is given by

Maximizing
the
posterior

$$\operatorname{argmax}_{\beta} p(\beta | \{t^{(i)}, x^{(i)}\}_{i=1}^N)$$

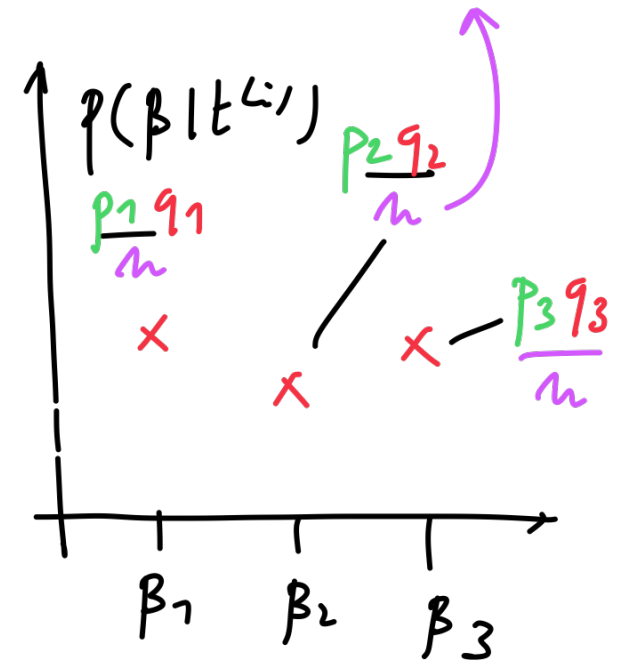
$$= \operatorname{argmax}_{\beta} p(\{t^{(i)}, x^{(i)}\}_{i=1}^N | \beta) p(\beta)$$

using the log to simplify the product
and the exponentials.

$$\operatorname{argmax}_{\beta} \log \left\{ \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2}{2\sigma^2}\right) \right\}$$

$$\times \prod_{j=1}^D \frac{1}{\sqrt{2\pi\lambda}} \exp\left(-\frac{\beta_j^2}{2\lambda^2}\right)$$

$$\begin{aligned} \mu &= p(t^{(i)} | \beta_1) p(\beta_1) \\ &+ p(t^{(i)} | \beta_2) p(\beta_2) \\ &+ p(t^{(i)} | \beta_3) p(\beta_3) \end{aligned}$$



$$\arg \max_{\beta} p(\{t^{(i)}, x^{(i)}\}_{i=1}^N | \beta) p(\beta)$$

$$= \arg \max_{\beta} \sum_{i=1}^N \left(\log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{1}{2\sigma^2} \left(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}) \right)^2 \right)$$

do not depend on β_j

$$\left\{ + \sum_{j=2}^D \left(\log \left(\frac{1}{\sqrt{2\pi}\lambda} \right) - \frac{1}{2\lambda^2} \beta_j^2 \right) \right\}$$

$$= \arg \max_{\beta} - \frac{1}{2\sigma^2} \sum_{i=1}^N \left(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}) \right)^2 - \frac{1}{2\lambda^2} \sum_{j=1}^D \beta_j^2 \quad *$$

Maximizing $-l(\beta)$ is equivalent to minimizing $l(\beta)$ from this

$$(*) = \underset{\beta}{\operatorname{arg\,min}} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2 + \frac{\sigma^2}{\lambda^2} \sum_{j=1}^D \beta_j^2 = \underline{\beta_{\text{MAP}}}$$

→ Under our assumptions on $\varepsilon^{(i)}$, β_j , finding β_{Ridge} is equivalent to finding the MAP estimator for β
 (Maximization of posterior distribution)

In short: MLE : β is found by maximizing the likelihood function

MAP : Some prior information on β (encoded in $p(\beta)$) is used to define a posterior distribution which is then maximized

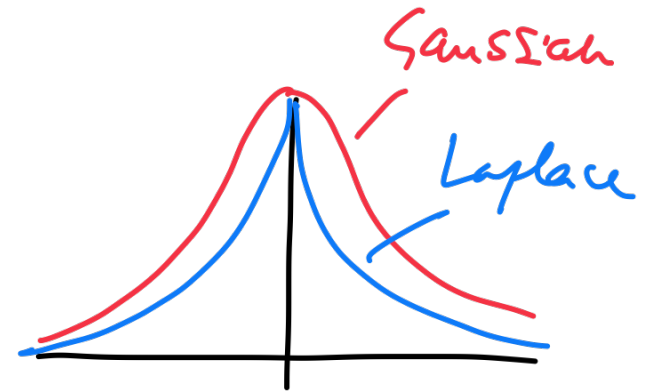
With respect to β .

→ The **Gaussian prior** is good at reducing the variance as it corresponds to the assumption that β_j are most likely small but we can do better.

→ In the Gaussian, there is not a big difference between being 0 and being very small.

To further emphasize the feature selection of the model (i.e. the fact that the β_j

should be exactly set to 0 whenever possible



We can use spikier distributions (e.g. Laplace)

→ Applying the previous (MAP) steps we get

$$\beta_{\text{MAP}} = \operatorname{argmax}_{\beta} \underbrace{p(\{t^{(i)}, x^{(i)}\}_{i=1}^N | \beta)}_{\text{likelihood}} p(\beta) \prod_{j=1}^D \frac{1}{2b} \exp\left(-\frac{|\beta_j|}{b}\right)$$

$$\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2}{2\sigma^2}\right)$$

$$\beta_{\text{MAP}} = \operatorname{argmin}_{\beta} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2 + \frac{\sigma^2}{b} \sum_{j=1}^D |\beta_j|$$

For a Laplace prior, the MAP estimator is the LASSO loss. minimizer of the

In short: under appropriate statistical assumptions,

β_{Ridge} = MAP with Gaussian prior

β_{LASSO} = MAP with Laplace prior

Geometric intuition

$$\text{Ridge } \ell(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2 + \lambda \sum_{j=1}^D \beta_j^2 \quad (*)$$

$$\text{LASSO } \ell(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2 + \lambda \sum_{j=1}^D |\beta_j|$$

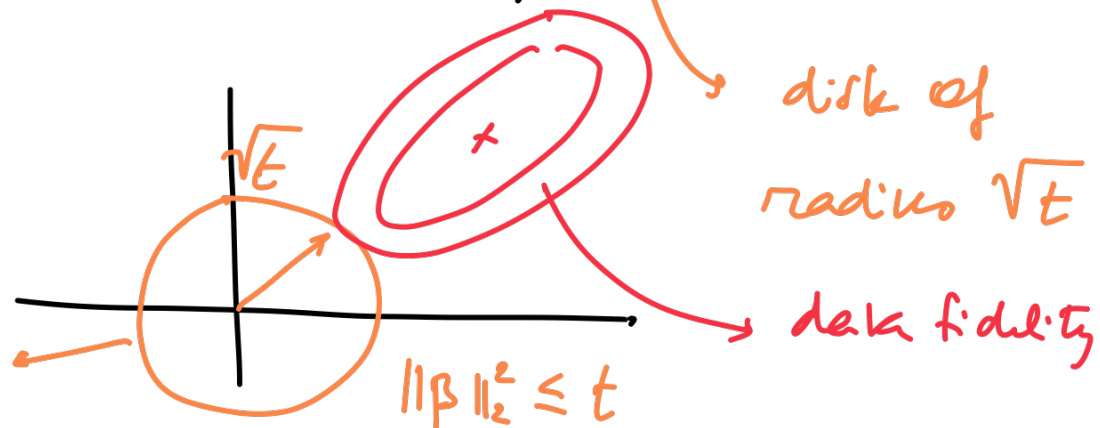
To get some additional intuition on how each formulation handles feature selection, let us consider the equivalent constrained formulations

$$\beta_{\text{Ridge}} = \min_{\beta} \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2$$

$$\text{s.t. } \sum_{j=1}^D |\beta_j|^2 \leq t \quad (**)$$

a larger value of t in **(**)** corresponds to a smaller λ in **(*)**

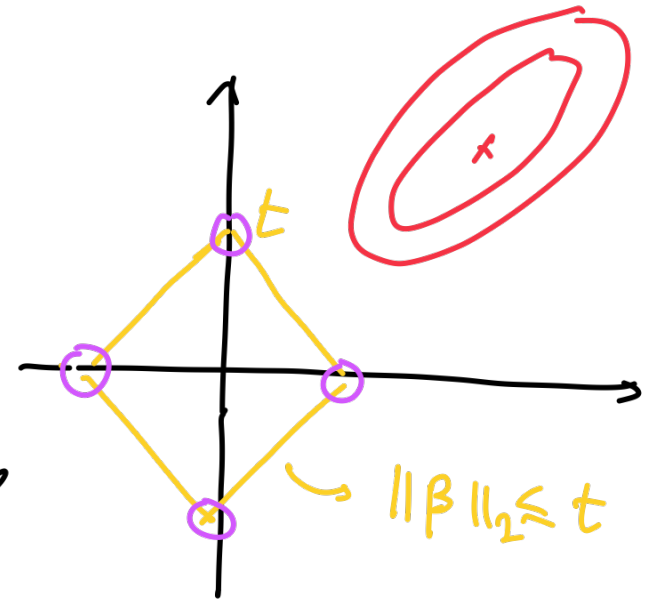
Ridge penalty



Such a constrained formulation can be equivalently defined for LASSO

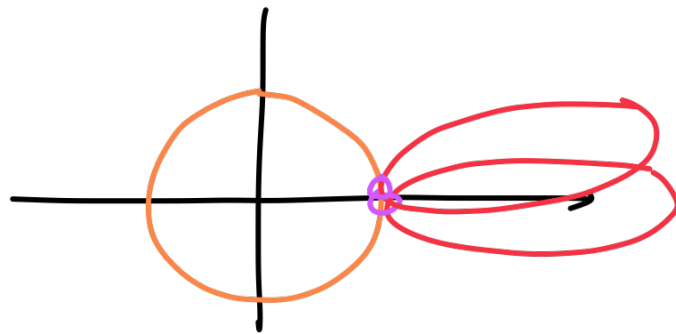
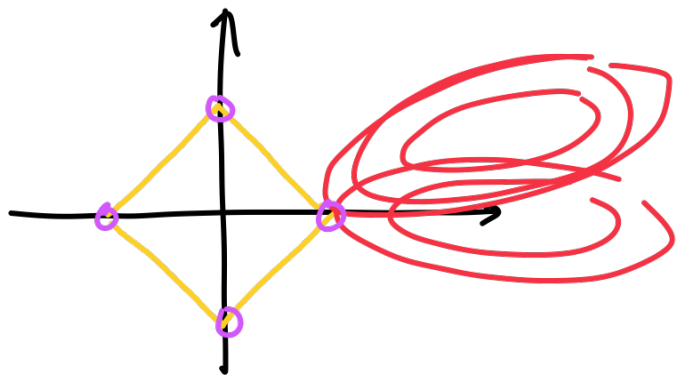
$$\beta_{\text{LASSO}} = \min \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2$$

$$\text{s.t. } \sum_{j=1}^D |\beta_j| \leq t$$



Both LASSO and Ridge can be understood as looking for the β that achieve the smallest possible value of the data fidelity term (objective corresponding to the red parabolas) within a certain ball (l_1 ball \rightarrow LASSO
 l_2 ball \rightarrow Ridge

→ Because the l_1 ball is more Spiky (such as Laplace vs Gauss)
the β inside the ball achieving the minimum of the
data fidelity term is more likely to be located on the
intersection with the axes than it is in the l_2 ball



Both Ridge and LASSO rely on a definition of the model complexity penalty based on the l_p norm for some p

In the case of Ridge we minimize the squared l_2 norm

$$\Omega(\beta) = \sum_{j=1}^D \beta_j^2 = \|\beta\|_2^2$$

In the case of LASSO we minimize the l_1 norm

$$\Omega(\beta) = \sum_{j=1}^D |\beta_j| = \|\beta\|_1$$

Of course other penalties can be defined for other values of p

$$\Omega(p) = \sum_{j=1}^D |\beta_j|^p = \|\beta\|_p^p = \left(\sqrt[p]{\sum_{j=1}^D |\beta_j|^p} \right)^p$$

