

Ridge Regression (a.k.a. Bayesian Linear Regression)

We have already seen:

$$y|x=x \sim N(\vec{w}^T \vec{x}, \sigma^2) \quad] \text{ Likelihood.}$$

Now let's put a prior: $\vec{w} \sim N(0, t^2 I) \quad] \text{ prior}$

$$w_i \sim N(0, t^2) \quad (\text{iid samples})$$

$$P(\vec{w}) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi t^2}} \exp \left[-\frac{(w_i - 0)^2}{2t^2} \right]$$

$$\begin{aligned} \hat{\vec{w}}_{MAP} &= \arg \max_{\vec{w}} \log P(\vec{w} | D) \\ &= \arg \max_{\vec{w}} \log \left[\frac{P(D|\vec{w})P(\vec{w})}{P(D)} \right] \\ &= \arg \max_{\vec{w}} \left[\log P(D|\vec{w}) + \log P(\vec{w}) \right] \\ &= \arg \max_{\vec{w}} \left[-\sum_{i=1}^n \frac{(y_i - \vec{w}^T \vec{x}_i)^2}{2\sigma^2} - \frac{n}{2} \log(2\pi\sigma^2) - \sum_{j=0}^d \frac{w_j^2}{2t^2} - \frac{1}{2} \log(2\pi t^2) \right] \\ &= \arg \min_{\vec{w}} \left[-\sum_{i=1}^n \frac{(y_i - \vec{w}^T \vec{x}_i)^2}{2\sigma^2} - \frac{n}{2} \log(2\pi\sigma^2) - \sum_{j=0}^d \frac{w_j^2}{2t^2} \right] \quad \text{const. wrt } w \\ &= \arg \min_{\vec{w}} \left[\sum_{i=1}^n (y_i - \vec{w}^T \vec{x}_i)^2 + \frac{2\sigma^2}{2t^2} \sum_{j=0}^d w_j^2 \right] \quad (\text{ans?}) \\ &= \arg \min_{\vec{w}} \left[\underbrace{\sum_{i=1}^n (y_i - \vec{w}^T \vec{x}_i)^2}_{\text{Average error}} + \underbrace{\frac{2\sigma^2}{2t^2} \sum_{j=0}^d w_j^2}_{\lambda \|\vec{w}\|_2^2} \right] \\ &= \arg \min_{\vec{w}} \left[\underbrace{\sum_{i=1}^n (y_i - \vec{w}^T \vec{x}_i)^2}_{\text{Avg. loss}} + \underbrace{\frac{\lambda}{n} \|\vec{w}\|_2^2}_{\text{regularization!}} \right] \\ &\boxed{\text{Loss} + \lambda \cdot \text{prior}} \end{aligned}$$

So, Bayesian loss = Avg. fitting error + $\lambda \cdot \text{Norm square of } \vec{w}$
 The norm of \vec{w} should be small (why?)

(12)

As $\lambda \rightarrow 0$, $t^2 \rightarrow \infty \Rightarrow$ Priors with a Gaussian having very large variance \rightarrow nearly uniform
 $\Rightarrow \hat{w}_{MAP} = \hat{w}_{MLE}$

As $\lambda \rightarrow \infty$, $\hat{w} \rightarrow 0$

$$\hat{w}^T x = w_0 + \underbrace{w_1 x_1 + w_2 x_2 + \dots + w_n x_n}_{\rightarrow 0}$$

In some sense, the regularizer encourages the optimization to pick a simple model - here pushes to mean.

With w_0 , we don't have any information about the data - we have only the y_i 's. The best that can be done is to take an average. So w_0 becomes the mean of y . So, \hat{y} becomes \bar{y}_{mean} .

Now let's get back to solving the minimization problem:

$$\hat{w}_{MAP} = \underset{w}{\operatorname{argmin}} \frac{1}{n} [(y - xw)^T (y - xw) + \lambda w^T w]$$

$$\begin{aligned} \frac{\partial}{\partial w} \left[\frac{1}{n} [(y - xw)^T (y - xw) + \lambda w^T w] \right] &= 0 \\ \Rightarrow \frac{\partial}{\partial w} \left[\frac{1}{n} \left(y^T y + w^T x^T x w - 2y^T x w + \lambda w^T w \right) \right] &= 0 \\ \Rightarrow 2w^T x^T x - 2y^T x + \cancel{2\lambda} w^T = 0 \end{aligned}$$

$$\Rightarrow x^T x = (x^T x + \cancel{\lambda I}) \cdot w$$

$$\Rightarrow \boxed{\hat{w}_{MAP} = (x^T x + \cancel{\lambda I})^{-1} x^T y}$$

Compare with $\hat{w}_{MLE} = (x^T x)^{-1} x^T y$

Modification of pseudo-inverse: $(x^T x)^{-1} \rightarrow (x^T x + \lambda I)^{-1}$

What adding an identity matrix does to matrix eigenvalues??

$$Ax = \lambda x$$

$$(A + I)x = Ax + Ix = (\lambda + 1)x$$

\hookrightarrow shifts eigenvalues up!!

So now the noninvertible problem becomes solvable.

Interpret the effect of regularizer.

Likelihood	Prior	Name
Gaussian	Uniform	Least Square
"	Gaussian	Ridge Regression
"	Laplace	Lasso
Laplace	Uniform	Robust Regression
Student	Uniform	Robust Regression.