

به نام خدا

علی قاسم زاده ۴۰۱۱۵۹۲۳۹

۱- داریس کر

$$f(y|x, \theta) = \mathcal{N}(y|x, \theta)$$

$$\Rightarrow y_i | x_i, \theta \sim \mathcal{N}(y_i|x_i, \theta)$$

$$\Rightarrow L(\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\theta}} \exp\left(-\frac{(y_i - y_i|x_i)^2}{2\theta}\right) \xrightarrow{\ln}$$

$$\ln L(\theta) = -\frac{1}{2} \left( \sum_{i=1}^n \ln(2\pi\theta) + \frac{(y_i - y_i|x_i)^2}{\theta} \right) \xrightarrow{\frac{\partial}{\partial \theta}}$$

$$-\frac{1}{2} \left( \sum_{i=1}^n \frac{1}{\theta} - \frac{(y_i - y_i|x_i)^2}{\theta^2} \right) = 0 \rightarrow$$

$$\frac{n}{\theta} - \frac{1}{\theta^2} \sum_{i=1}^n \frac{(y_i - y_i|x_i)^2}{\theta} = 0 \rightarrow n\theta = \sum_{i=1}^n \frac{(y_i - y_i|x_i)^2}{\theta}$$

$$\rightarrow \boxed{\hat{\theta}_{ML} = \frac{1}{n} \sum_{i=1}^n \frac{(y_i - y_i|x_i)^2}{\theta}}$$

$$y = f(m) + \varepsilon, \varepsilon \sim \mathcal{N}(0, \sigma^2) \quad 1-2 \text{ داریس کر}$$

$$\text{Err}(m) = \mathbb{E}_{D, \varepsilon} \{ (y - \hat{f}(m))^2 \} = \mathbb{E}_{D, \varepsilon} \{ (f(m) + \varepsilon - \hat{f}(m))^2 \} =$$

$$\underbrace{\mathbb{E}_{D, \varepsilon} \{ \varepsilon^2 \}}_{\sigma^2} + \underbrace{\mathbb{E}_{D, \varepsilon} \{ (f(m) - \mathbb{E}\{\hat{f}(m)\})^2 \}}_{\text{Bias}^2(\hat{f}(m))} + \underbrace{\mathbb{E}_{D, \varepsilon} \{ (\hat{f}(m) - \mathbb{E}\{\hat{f}(m)\})^2 \}}_{\text{Var}_D(\hat{f}(m))}$$

همچنین عبارت بالا یعنی ترم حا سادۀ تر نیز محاسبه شود، بزر مثال  $E\{\varepsilon^2\}_{D, \varepsilon}$  همان  $E\{\varepsilon^2\}$  است و ...

$$\Rightarrow \text{Err}(x) = \text{Bias}(\hat{f}(x))^2 + \text{Var}(\hat{f}(x)) + \sigma^2$$

۲-۲. حال با سیمپل

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

$$D = \{(x_i, y_i), i=1, \dots, n\}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = 0 \quad \frac{1}{n} \sum_{i=1}^n x_i^2 = S_x^2$$

$$\hat{f}_1(x) = \frac{1}{n} \sum_{i=1}^n y_i \quad \hat{f}_r(x) = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$E\{\hat{f}_1(x)\} = E\left\{\frac{1}{n} \sum_{i=1}^n y_i\right\} = \beta_0 + \beta_1 \bar{x} = \beta_0$$

$$\Rightarrow \text{Bias}(\hat{f}_1(x)) = \beta_0 + \beta_1 x - \beta_0 = \beta_1 x$$

$$\text{Var}(\hat{f}_1(x)) = \text{Var}(\bar{y}) = \frac{\sigma^2}{n} \rightarrow \text{Err}(x) = \beta_1^2 x^2 + \frac{\sigma^2}{n} + \sigma^2$$

بزر  $\hat{f}_r(x)$  باز سیمپل

$$\hat{f}_r(x) = \hat{\beta}_0 + \hat{\beta}_1 x, (\hat{\beta}_0, \hat{\beta}_1) = \underset{\beta_0, \beta_1}{\text{argmin}} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad \bar{x} = 0$$

$$\leadsto E\{\hat{\beta}_0\} = \beta_0, E\{\hat{\beta}_1\} = \beta_1, \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = 0$$

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2}{n}, \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{n S_x^2}$$



$$\Rightarrow \text{Err}(n) = ? \quad \text{B } E\{\hat{f}_r(n)\} = \beta_0 + \beta_1 x = f(n) \rightarrow \text{Bias}$$

$$\Rightarrow \text{Err}(n) = 0 + \left( \frac{\sigma^2}{n} + n^2 \frac{\sigma^2}{n s_x^2} \right) + \sigma^2$$

←  $\hat{f}_1$  مدل ساده‌تری است ولی بایاس زیادی دارد ولی در عوض واریانس کمتری  $\left(\frac{\sigma^2}{n}\right)$  دارد.  
 در مدل  $\hat{f}_r$  پیچیده‌تر است و بایاس ندارد ولی در عوض واریانس بیشتری دارد که نسبت دارد.

$$(\hat{\beta}_0^\lambda, \hat{\beta}_1^\lambda) = \underset{\beta_0, \beta_1}{\text{argmin}} \sum_i (y_i - \beta_0 - \beta_1 x_i)^2 + \lambda \beta_1^2 \rightarrow \text{Ridge Regression}^{2-2}$$

؛ مشتق‌گیری و برابر صفر گذاشتن داریم که

$$\hat{\beta}_0^\lambda = \bar{y}, \quad \hat{\beta}_1^\lambda = \frac{n s_x^2}{n s_x^2 + \lambda} \hat{\beta}_1^{\text{ols}}$$

که  $\hat{\beta}_1^{\text{ols}}$  همان  $\beta_1$  در سری دوم یعنی دوم این‌ها است.

$$\Rightarrow E\{\hat{f}_\lambda(n)\} = \beta_0 + \frac{n s_x^2}{n s_x^2 + \lambda} \beta_1 x \rightarrow \text{bias} = \beta_1^2 x^2 \frac{\lambda}{(n s_x^2 + \lambda)^2}$$

$$\text{Var}(\hat{f}_\lambda(n)) = \frac{\sigma^2}{n} + n^2 \left( \frac{\sigma^2 n s_x^2}{(n s_x^2 + \lambda)^2} \right)$$

حالا اگر  $\lambda \rightarrow 0$  ، بایاس صفر میشه ، ridge  $\rightarrow$  ols

اگر  $\lambda \rightarrow \infty$  ، بایاس ~~بیشتر~~ ثابت میشه و واریانس میشه  $\frac{\sigma^2}{n}$

← زیاد کردن  $\lambda$  بایاس را بیشتر میکنه ولی واریانس را کمتر می‌کنه.

$$y_i = \sin(\pi x_i) + \varepsilon_i, \quad x_i \sim \text{unif}([0, 1]) \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n y_i \Rightarrow E\{\hat{f}(x)\} = \int_0^1 \sin(\pi x) dx = \frac{2}{\pi}$$

$$\Rightarrow \mu = \frac{2}{\pi}, \quad V_f = \text{Var}(\sin(\pi x)) = \int_0^1 \sin^2(\pi x) dx - \mu^2 = \frac{1}{2} - \frac{4}{\pi^2}$$

$$\text{Var}(\bar{y}) = \frac{1}{n} (V_f + \sigma^2)$$

$$\text{Gen Err} = E_{x, D, \varepsilon} \left\{ (y - \bar{y})^2 \right\} = E \left\{ (\sin(\pi x) - \mu)^2 \right\} + \text{Var}(\bar{y}) + \sigma^2$$

$$\Rightarrow \text{Gen Err} = \underbrace{\left( \frac{1}{2} - \frac{4}{\pi^2} \right)}_{\text{misspecification bias}} + \underbrace{\frac{1}{n} \left( \frac{1}{2} - \frac{4}{\pi^2} + \sigma^2 \right)}_{\text{est var}} + \underbrace{\sigma^2}_{\text{irreducible noise}}$$

وقتی  $n \rightarrow \infty$  ترم  $\frac{1}{n}$  در ناپدید می شود و اور ما برابر خواهد بود با  $\frac{1}{2} - \frac{4}{\pi^2} + \sigma^2$

$$\frac{\partial \ell}{\partial c} = \frac{\partial \ell}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial c} = -(y - \hat{y}) \cdot 1 = v^T h + c - y$$

$$\frac{\partial \ell}{\partial v} = \frac{\partial \ell}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial v} = -(y - \hat{y}) h \in \mathbb{R}^m$$

حالا بایراد  $h$  مشتق بگیریم خواصی داریم

$$\frac{\partial \ell}{\partial h} = \frac{\partial \ell}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial h} = -(y - \hat{y}) v \in \mathbb{R}^m$$



بیل Relu داریم.

$$\frac{\partial \ell}{\partial h} = \frac{\partial \ell}{\partial z} \cdot \frac{\partial z}{\partial h} = -(y - \hat{y}) \cdot \nabla \in \mathbb{R}^m$$

$$\sigma'(z_j) = \begin{cases} 1 & z_j > 0 \\ 0 & z_j < 0 \end{cases}$$

$$\Rightarrow \frac{\partial \ell}{\partial z} = \frac{\partial \ell}{\partial h} \odot \sigma'(z) = -(y - \hat{y}) \cdot \nabla \odot 1\{Ux + b > 0\}$$

$$\frac{\partial \ell}{\partial b} = \frac{\partial \ell}{\partial z} = -(y - \hat{y}) \cdot \nabla \odot 1\{Ux + b > 0\}$$

$$\frac{\partial \ell}{\partial U} = \frac{\partial \ell}{\partial z} \cdot \frac{\partial z}{\partial U} = -(y - \hat{y}) \cdot \nabla \odot 1\{Ux + b > 0\} X^T$$

$m \times d$  matrix

تدریجاً مقادیر بی‌بسته، احتمال اینکه  $z_j$  صفر باشد صفر.

از طرفی می‌توانیم هر subgradient نیز استفاده کنیم.

از طرفی Relu جابر negative را دیگر آپدیت نمی‌کند و dead می‌شوند.

دباعث sparsity بیشتر، generalization بهتر در  $\mathbb{R}^n$  و یک سری از کسرها می‌شود.

۴- منظم‌سازی، پنهان در روش گاهش گزینان.

۱.۴ جواب mle عبارت است از:

$$\nabla_{\beta} (y - X\beta)^T (y - X\beta) = 0 \rightarrow X^T X \beta = X^T y \rightarrow$$

$$\beta_{\min} = (X^T X)^{-1} X^T y \rightarrow$$

با فرضی بودن  $X^T X$

حالا در روش گاهش گزینان داریم:

$$\beta^{(t+1)} = \beta^{(t)} - \eta \nabla \|y - X\beta^{(t)}\|^2 = \beta^{(t)} - 2\eta X^T (X\beta^{(t)} - y) = (I - 2\eta X^T X) \beta^{(t)} + 2\eta X^T y$$

$$\beta^{(t)} = (I - \eta A)^t \beta^{(0)} + \eta \sum_{k=0}^{t-1} (I - \eta A)^k X^T Y$$

$\downarrow$   $X^T X$        $\downarrow$   $X^T X$

آنکه  $\eta$  را به اندازه کافی کوچک بگیریم که تمامی معاینه‌ها در  $I - \eta A$  بین  $(1, -1)$  بقیه.

$$\Rightarrow (I - \eta A)^t \rightarrow 0 \text{ when } t \rightarrow \infty \Rightarrow$$

$$M^t \rightarrow V(I - \eta \Lambda)V^T, \quad (1 - \eta \lambda_j)^t \xrightarrow{t \rightarrow \infty} 0$$

$$\Rightarrow M^t = (I - \eta A)^t = V(I - \eta \Lambda)^t V^T \rightarrow 0$$

$$\sum_{k=0}^{\infty} M^k = (I - M)^{-1} = (I - (I - \eta A))^{-1} = (\eta A)^{-1}$$

$$\Rightarrow \lim_{t \rightarrow \infty} \beta^{(t)} = \eta (\eta A)^{-1} X^T Y = A^{-1} X^T Y = (X^T X)^{-1} X^T Y$$

← به همان جواب  $\argmin$  می‌رسیم.

$$\beta^{(t)} = (I - (I - \eta A)^t) A^{-1} X^T Y$$

۲-۴.  $\lambda_j$  را

$$\hat{\beta}_{\text{ridge}}(\lambda) = (A + \lambda I)^{-1} X^T Y, \quad A = V \Lambda V^T$$

$$\beta_j^{(t)} = \left( 1 - (1 - \eta \lambda_j)^t \right) \frac{V_j^T X^T Y}{\lambda_j} \Rightarrow$$

$$\beta_j^{\text{ridge}} = \frac{1}{\lambda + \lambda_j} (V_j^T X^T Y)$$

کوچک‌ترین



حالا برای اینکه این دو معادله بشوند، باید داشته باشیم:

$$\frac{(1 - (1 - 2\eta \lambda_j)^t)}{\lambda_j} = \frac{1}{\lambda + \lambda_j} \rightarrow \lambda = \frac{\lambda_j (1 - 2\eta \lambda_j)^t}{1 - (1 - 2\eta \lambda_j)^t}$$

برای هر  $\lambda_j$  ای.

متوقف شدن در کام  $t$ ام معادل است با ridge reg که  $\lambda$  آن مقدر بالا را داشته باشد.

۳-۴. مثل مدل خطی روی داده‌های بی‌س روی شبکه‌های عصبی علاوه بر کمینه کردن خطا

با تنظیماتی مثل تعداد گام بروزرسانی، نرخ یادگیری، batch-size به صورت غیر مستقیم

پارامترها را به سمتی می‌برند که تابع ساده‌تر بیابند، ~~weight~~

از طریق  $\text{early stopping}$  داشته باشیم و آموزش را قبل از همگرایی متوقف کنیم، از بروز  $\text{overfit}$  جلوگیری می‌شود. کار ridge regression برابر ما می‌کند.

۵ - داریم که

$$\hat{\theta} = (1 - \frac{c}{\|y\|^2}) y \rightarrow \hat{\theta} - \theta = (y - \theta) - \frac{c}{\|y\|^2} y \Rightarrow$$

$$\|\hat{\theta} - \theta\|^2 = \|y - \theta\|^2 + c^2 \frac{\|y\|^2}{\|y\|^4} - 2c \frac{(y - \theta)^T y}{\|y\|^2} \xrightarrow{E}$$

$$E\{\|\hat{\theta} - \theta\|^2\} = E\{\|y - \theta\|^2\} = d$$

stein identity

$$E\{(y - \theta)^T f(y)\} = E\{\text{div } f(y)\}$$

$$f(y) = \frac{y}{\|y\|^2} \rightarrow \text{div } \frac{y}{\|y\|^2} = \sum_{i=1}^d \frac{\partial}{\partial y_i} \left( \frac{y_i}{\|y\|^2} \right) = \frac{d-2}{\|y\|^2}$$

$$\Rightarrow E\left\{\frac{(y-\theta)^T y}{\|y\|^2}\right\} = E\left\{\frac{d-2}{\|y\|^2}\right\}$$

$$\Rightarrow E\{\|\hat{\theta} - \theta\|^2\} = d - 2c E\left\{\frac{d-2}{\|y\|^2}\right\} + c^2 E\left\{\frac{1}{\|y\|^2}\right\} =$$

$$d + (c^2 - 2c(d-2)) E\left\{\frac{1}{\|y\|^2}\right\}$$

۲-۵. حالا مقدار بهینه  $c$  را بخواهیم پیدا کنیم.

$$2c - 2(d-2) = 0 \rightarrow c^* = d-2$$

$$\Rightarrow d - (d-2)^2 E\left\{\frac{1}{\|y\|^2}\right\}$$

برای  $mle$  که بخواهیم بنویسیم درمیان که  $\hat{\theta}_{mle} = y$  حالا به بلوریه کسی  
نقطه از اتفاق افتاد که مرکز کسی باشه بیشترین احتمال را داره ✓

$$\Rightarrow \text{risk}, E\{\|y - \theta\|^2\} = d$$

ما داریم که  $c = d-2$

$$\rightarrow d - (d-2)^2 E\left\{\frac{1}{\|y\|^2}\right\} < d$$

۳-۵. برای  $d \geq 3$ ، James-stein مقار  $\hat{\theta}_{JS}$  ریسک کمتر از  $\hat{\theta}_{mle}$  می‌کند.

$$\hat{\theta}_{JS} = \left(1 - \frac{d-2}{\|y\|^2}\right) y \rightarrow$$

ریسک کمتر داره.

این از  $mle$  در  $mean-squared error$  بهتره.