

SDS 385: Stat Models for Big Data Lecture 2: Linear Regression

Purnamrita Sarkar Department of Statistics and Data Science The University of Texas at Austin

https://psarkar.github.io/teaching

Technical things you need for this lecture

- What is the log likelihood and what is MLE– and how to get the MLE?
- How to do derivatives w.r.t a vector. (See section 2.4) https: //www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf
- What is the Hessian of a function?
- Eigenvalues of a matrix and how they change when you add a multiple of identity.
- What is a Positive Semi-Definite (PSD) matrix?
- What is a convex function (Jensen's inequality) and what operations preserves convexity? See 2.3 in

https://web.stanford.edu/~boyd/cvxbook/bv_cvxbook.pdf

- What is a strongly convex function and how can you relate to the Hessian.
- What is the mean value theorem in Calculus, we will do the vector version of this.
- What is a Lipschitz continuous function?

Linear regression: recap

Given *n* pairs $(\mathbf{x}_i, y_i) \in \Re^{p+1 \times 1}$, consider the model:

$$\mathbf{y} = \mathbf{X}\mathbf{\beta} + \boldsymbol{\epsilon} \qquad \epsilon_i \sim N(0, \sigma^2)$$

where:

$$\mathbf{y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}, \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}, \text{ and } \mathbf{x} = \begin{bmatrix} 1 & x_{12} & \dots & x_{1p} \\ 1 & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n2} & \dots & x_{np} \end{bmatrix}$$

• X, y are given, you need to estimate β .

MLE - recap

$$f(\mathbf{y}|\mathbf{X}; \mathbf{eta}) \propto \exp(rac{-(\mathbf{y} - \mathbf{X}\mathbf{eta})^T(\mathbf{y} - \mathbf{X}\mathbf{eta})}{2\sigma^2})$$

• Take Log, we can get:

$$\frac{-(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^{\mathsf{T}}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2}$$
(1)

• Same drill - differentiate and set it to zero.

$$-\boldsymbol{X}^{T}(\boldsymbol{y}-\boldsymbol{X}\hat{\boldsymbol{\beta}})=0\rightarrow\boldsymbol{X}^{T}\boldsymbol{X}\hat{\boldsymbol{\beta}}=\boldsymbol{X}^{T}\boldsymbol{y}\rightarrow\hat{\boldsymbol{\beta}}=(\boldsymbol{X}^{T}\boldsymbol{X})^{-1}\boldsymbol{X}^{T}\boldsymbol{y}$$

• What happens when $p \gg n$? $\mathbf{X}^T \mathbf{X}$ is not invertible.

- Add a prior to β , i.e. $\beta \sim N(0, \lambda I_p)$, or think of adding a regularization that penalizes large values of $\beta^T \beta$.
- So now we have:

$$f(\mathbf{y}|\mathbf{X}, \mathbf{\beta}) \propto \exp(\frac{-(\mathbf{y} - \mathbf{X}\mathbf{\beta})^{\mathsf{T}}(\mathbf{y} - \mathbf{X}\mathbf{\beta})}{2\sigma^2} - \lambda \mathbf{\beta}^{\mathsf{T}}\mathbf{\beta})$$

• Differentiating and setting to zero gives:

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{X}^{T}\boldsymbol{X} + \lambda \boldsymbol{I}_{p})^{-1}\boldsymbol{X}^{T}\boldsymbol{y}$$

• Phew! – no issues with invertibility of $X^T X$

Exact computation

- If **X** was dense, how much time would the computation of **X**^T**X** take?
- Wait, what is dense?
- Well dense means, **X** has about $\Theta(np)$ non-zero elements.



Figure 1: Dense matrix multiplication¹

• So O(n) computation for each of p^2 entries, and hence np^2 .

¹Borrowed from Cho-Jui Hsieh's classnotes at UC-Davis.

Sparse matrix data structures

- How do you store a sparse vector?
- All you need is two vectors: one is of the indices of nonzero elements and one is the values.



 $C_{ij} = \sum_{k} A_{ik} \cdot B_{kj}$

Figure 2: Sparse matrix multiplication²

• If A has nnz non-zeroes, then worst case, the complexity is $O(n \times nnz)$ operations for multiplying a sparse matrix with another dense matrix.

²Borrowed from Grey Ballard and Alex Druinsky, SIAM conf. on Lin. Algenbra

- Inverting a $p \times p$ matrix takes $O(p^3)$ time.
- Alternatives: use linear solvers of the form Au = v.
- Here $A = \mathbf{X}^T \mathbf{X} + \lambda I_p$, $\mathbf{v} = \mathbf{X}^T \mathbf{y}$ and $\mathbf{u} = \beta$.
- Unless your matrix A has some structure, linear solvers can also be expensive. However, if it does have structure, e.g. its diagonally dominant, etc, then there are nearly linear time solvers.
- Typically for regression, we don't expect to have such structure.
- So, what can be done?

• Lets talk about gradient descent type methods.

• Model:
$$\sum_{i=1}^{n} f(\underbrace{x_i}_{data}; \underbrace{\beta}_{parameter})$$

• Example of *f*: negative log-likelihood over iid data-points, e.g. linear regression, logistic regression, etc.

• Goal:
$$\hat{\beta} = \arg\min_{\beta} f(x_i; \beta)$$

• Lets deal with convex loss functions.

Convex functions



Figure 3: A convex function

 $\forall \alpha \in [0,1], f(\alpha x + (1-\alpha)y) \le \alpha f(x) + (1-\alpha)f(y)$

$$f(\alpha x + (1 - \alpha)y) = (\alpha x + (1 - \alpha)y)^{2}$$

= $\alpha^{2}x^{2} + (1 - \alpha)^{2}y^{2} + 2\alpha(1 - \alpha)xy$
 $\leq \alpha^{2}x^{2} + (1 - \alpha)^{2}y^{2} + \alpha(1 - \alpha)(x^{2} + y^{2})$
= $\alpha x^{2} + (1 - \alpha)y^{2}$

• Where did I use $\alpha \in [0, 1]$?

Convex functions f(y) = |y|



Figure 4: f(y) = |y|

$$f(\alpha x + (1 - \alpha)y) = |\alpha x + (1 - \alpha)y|$$
$$\leq |\alpha x| + |(1 - \alpha)y|$$
$$\leq \alpha |x| + (1 - \alpha)|y$$

Theorem

Consider an optimization problem $\min_{x} f(x)$ where f is convex. Let x^* be a local minima. Prove that it is also a global minima.

Proof.

- By definition, $\exists p > 0$, such that $\forall x \in B(x^*, p), f(x) \ge f(x^*)$.
- If x^* is not the global optima, $z \notin B(x^*, p)$ such that $f(z) < f(x^*)$.
- Take $t \in [0, 1]$ and the point $y = tx^* + (1 t)z$. $f(y) \le tf(x^*) + (1 - t)f(z) < f(x^*)$
- Now $|y x^*| = (1 t)|z x^*|$. If we take t large enough such that $(1 t)|z x^*| \le p$, then $y \in B(x^*, p)$ but $f(y) < f(x^*)$, which is a contradiction.

- Non-negative linear combinations of convex functions is also convex.
 - For example af(x) + bg(x) where f, g are both convex and a, b ≥ 0 is also convex.
- A convex function composed with an affine function is also convex.
- Point-wise maxima of convex functions is convex.

- Compositions of convex functions not necessarily convex
- f,g convex.
 - Is f g convex?
 - Is fg convex?

• First order:

$$\langle x - y, \nabla f(x) - \nabla f(y) \rangle \geq 0$$

• Second order:

 $\nabla^2 f(x) \succeq 0$

- Example: $f(x) = x^2$. $(x y)^2 > 0$ and $f''(x) = 2 \ge 0$.
- Example: $f(x) = \log(1 + e^x)$.
 - f'(x) = 1/(1 + e^{-x}) is monotonically increasing with x and so the first order condition is satisfied.
 - Second order: $f''(x) = f(x)(1 f(x)) \ge 0$

• First order:

$$\langle x - y, \nabla f(x) - \nabla f(y) \rangle \ge \mu ||x - y||^2$$

• Second order:

$$\nabla^2 f(x) \succeq \mu f$$

• So you add a margin to each inequality.

$$\beta \leftarrow \beta - \alpha \nabla f(\beta)$$



Figure 5: Convex function minimization with gradient descent ¹⁷

Step size



Figure 6: Choice of step size is crucial

Reality of gradient descent for a nonconvex function



Figure 7: Very nonconvex function



Figure 8: Newton Raphson⁴

⁴Borrowed from Nick Alger, math.stackexchange.com

Newton Raphson cont.

- GD takes into account only first order information.
- NR also takes second order information.
- In particular it uses the Hessian,

$$H(i,j) = \frac{\partial^2 f}{\partial \theta_i \partial \theta_j}$$
, where $i, j \in \{1, \dots, k\}$,

• Lets try to minimize a quadratic function:

$$f = \mathbf{a} + \mathbf{b}^T \boldsymbol{\theta} + \frac{1}{2} \boldsymbol{\theta}^T C \boldsymbol{\theta}.$$

- C is positive semidefinite and so this is a convex function.
- We can minimize the function by differentiating it and by setting the result equal to 0:

$$abla f(oldsymbol{ heta}^*) = \mathbf{b} + Coldsymbol{ heta}^* = 0$$

 $oldsymbol{ heta}^* = -C^{-1}\mathbf{b}$

• In the neighborhood of θ_t , we can use the approximation:

$$f(\boldsymbol{\theta}^{(t)} + \mathbf{h}) \approx f(\boldsymbol{\theta}^{(t)}) + \nabla f(\boldsymbol{\theta}^{(t)})^{T} \mathbf{h} + \frac{1}{2} \mathbf{h}^{T} H(\boldsymbol{\theta}^{(t)}) \mathbf{h}.$$
(2)

• Therefore the general updating rule is

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - H^{-1}(\boldsymbol{\theta}^{(t)}) \cdot \nabla f(\boldsymbol{\theta}^{(t)})$$

• You can use a stepsize here as well.

• for t = 1 : T (or until convergence)

• Do
$$\beta_{t+1} \leftarrow \beta_t - \alpha \nabla f(\beta)$$

Theorem

Let β^* is the global minima, and the second derivative is bounded as $\mu I \leq H(\beta) \leq LI$. Then with $\alpha = 2/(L + \mu)$, gradient descent converges geometrically, i.e.

$$\|\beta_{t+1} - \beta^*\| \le \frac{L-\mu}{L+\mu} \|\beta_t - \beta^*\|$$

Proof

• Lets look at the distance from the optima:

$$\beta_{t+1} - \beta^* = \beta_t - \beta^* - \alpha(\nabla f(\beta_t) - \nabla f(\beta^*))$$
$$= \beta_t - \beta^* - \alpha H(z_t)(\beta_t - \beta^*)$$
$$= (I - \alpha \nabla^2 f(z_t))(\beta_t - \beta^*)$$

• Now take norm of both sides and use Triangle.

$$\begin{aligned} \|\beta_{t+1} - \beta^*\| &\leq \|I - \alpha H(z_t)\| \|\beta_t - \beta^*\| \\ &\leq \max(|1 - \alpha \mu|, |1 - \alpha L|) \|\beta_t - \beta^*\| \end{aligned}$$

Proof



• Pick
$$\alpha = \arg \min \max(|1 - \alpha \mu|, |1 - \alpha L|) = 2/(L + \mu)$$

• Set $\alpha = 2/(L + \mu)$. You get

$$\|\beta_{t+1} - \beta^*\| \le \frac{L-\mu}{L+\mu} \|\beta_t - \beta^*\|$$

• Finally after T iterations, we have:

$$\|\beta_{T+1} - \beta^*\| \le \left(\frac{L-\mu}{L+\mu}\right)^T \|\beta_1 - \beta^*\|$$

• This is a typical "linear" contraction result.

Theorem

Let β^* is the global minima, and the second derivative is L Lipschitz, i.e. $\|H(x) - H(x')\| \le \kappa \|x - x'\|$ and $\|H^{-1}\| \le 1/\mu$. Then with $\alpha = 1$, Newton Raphson converges quadratically, i.e.

$$\|\beta_{t+1} - \beta^*\| \le \kappa/\mu \|\beta_t - \beta^*\|^2$$

• Note that this is useful only when $\|eta_{t+1} - eta^*\| \ll 1$

Proof

$$\beta_{t+1} - \beta^* = \beta_t - \beta^* - H^{-1}(\beta_t)(\nabla f(\beta_t) - \nabla f(\beta^*)) \\ = \beta_t - \beta^* - H^{-1}(\beta_t)H(z_t)(\beta_t - \beta^*) \\ = (I - H^{-1}(\beta_t)H(z_t))(\beta_t - \beta^*) \\ = H^{-1}(\beta_t)(H(\beta_t) - H(z_t))(\beta_t - \beta^*) \\ \|\beta_{t+1} - \beta^*\| \le \|H^{-1}(\beta_t)\|\|H(\beta_t) - H(z_t)\|\|\beta_t - \beta^*\| \\ \le \kappa/\mu\|\beta_t - z_t\|\|\beta_t - \beta^*\| \\ \le \kappa/\mu\|\beta_t - \beta^*\|^2$$

- You have to calculate the gradient every iteration.
- Take ridge regression.
- You want to minimize $1/n\left((\boldsymbol{y} \boldsymbol{X}\boldsymbol{\beta})^{T}(\boldsymbol{y} \boldsymbol{X}\boldsymbol{\beta}) \lambda \boldsymbol{\beta}^{T}\boldsymbol{\beta}\right)$
- Take a derivative: $(-2\boldsymbol{X}^{T}(\boldsymbol{y}-\boldsymbol{X}\boldsymbol{\beta})-2\lambda\boldsymbol{\beta})/n$
- Grad descent update takes $\beta_{t+1} \leftarrow \beta_t + \alpha(\boldsymbol{X}^T(\boldsymbol{y} \boldsymbol{X}\beta_t) + \lambda\beta_t)$
- What is the complexity?
 - Trick: first compute $\boldsymbol{y} \boldsymbol{X}\boldsymbol{\beta}$.
 - np for matrix vector multiplication, nnz(X) for sparse matrix vector multiplication.
 - Remember the examples with humongous *n* and *p*?

Cho-Jui Hsieh and Christopher De Sa's large scale ML classes.