File   Edit   View   Run   Kernel   Nbgrader   Settings   Help

Markdown ∨

# Final exam review

What did we learn in this course?

1. Python basics:

   - Python_1.ipynb
   - Python_2.ipynb
2. Numpy basics:

   - NumpyTutorial.ipynb
3. Linear Regression by vector derivatives:

   - LinearModels.ipynb,
   - PlaneFitProblem,
   - Hessians,
   - Practice Problems for Midterm 1,
4. Optimization by Gradient descent:

   - ContinuousOptimization.ipynb
5. 1-Layer Neural Network: Perceptron3.ipynb

6. Pre-midterm review: Practice Problems

minimize
Quadratic form
of vectors
Any _vector_ expression
Quadratic
take its derivative
and equate it to zero

## Decision Theory:

$$\frac{\partial}{\partial \underline{x}}\left[ \underbrace{\underline{x}^T Q \underline{x}}_{} + \underline{b}^T \underline{x} + c \right.$$

$$\underline{x}^T(Q + Q^T) + \underline{b}^T + 0 = 0$$

$$\underline{x}^T = -\underline{b}^T (Q + Q^T)^{-1}$$

$$\underline{x} = -(Q + Q^T)^{-T} \underline{b}$$

- Lecture notes: decision-theory.pdf. Additional resources:
- Chapter 2 of MLStory book

## Bayes Rule

Question:

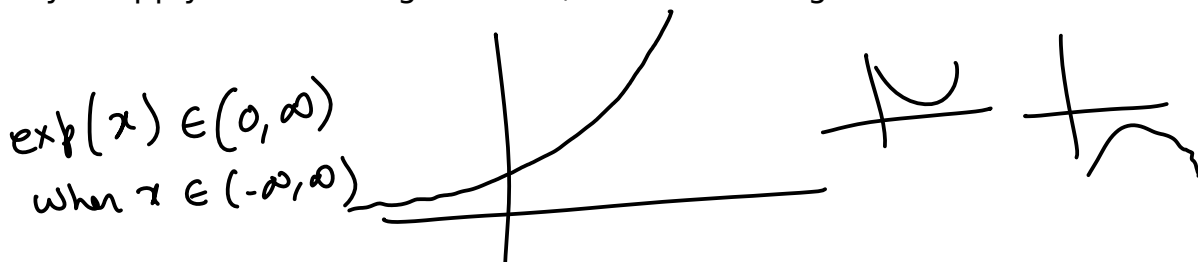Given two random variables $X$ and $Y$ specify the relationship between $P(X|Y)$ and $P(Y$

Answer

$P(X|Y)P(Y)$

Regularization can be interpretted as application for Bayes theorem $\begin{pmatrix} \text{Maximum-a} \\ \text{-posterior} \\ \text{estimate} \end{pmatrix}$

MAP estimate

$$\underset{w}{\arg\min} \quad L(D; w) + \lambda \|w\|_2^2 \quad \Longleftrightarrow \quad \underset{w}{\arg\max} \; P(W|D)$$

$$\in [0,1]$$

$$\underset{w}{\arg\max} \quad -L(D; w) - \lambda \|w\|_2^2$$

$$\in (-\infty, \infty)$$

In optimization, the optimal value/argument stays unchanged if you apply an increasing function. If you apply a decreasing function, the max changes to min and min changes to max.
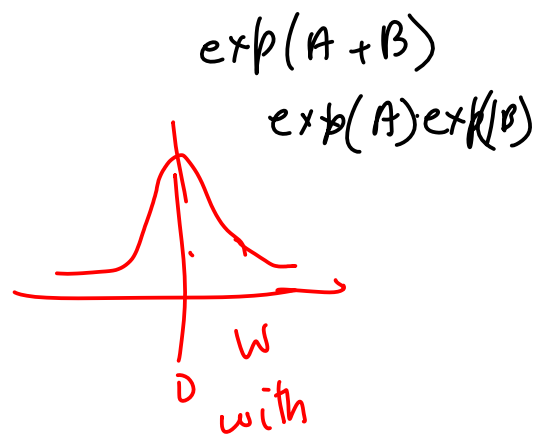
$$\exp(x) \in (0, \infty)$$
$$\text{when } x \in (-\infty, \infty)$$

$$\underset{w}{\arg\max} \; \frac{1}{Z} \exp\left(-L(D; w) - \lambda \|w\|_2^2\right) \quad \in [0, \infty)$$

$P(W|D)$ posterior

where $Z$ is a normalization factor so that

$$\sum_w \exp(\dots) = Z$$

Gaussian dist

$$\underset{w}{\arg\max} \; \frac{1}{Z} \underbrace{\exp(-L(D; w))}_{P(D|W)} \underbrace{\exp\left(-\|w\|^2 \cdot \frac{1}{\lambda}\right)}_{P(W)}$$

$$\exp(A+B)$$
$$\exp(A)\exp(B)$$

Evidence

likelihood       Prior

$D$ with

$$STD \cdot \frac{1}{\lambda} = 2\sigma^2$$

$$\Rightarrow \sigma = \frac{1}{2\sqrt{\lambda}}$$

# Vector Jacobian product

Reverse mode differentiation
And it assumes the final
output of the computation
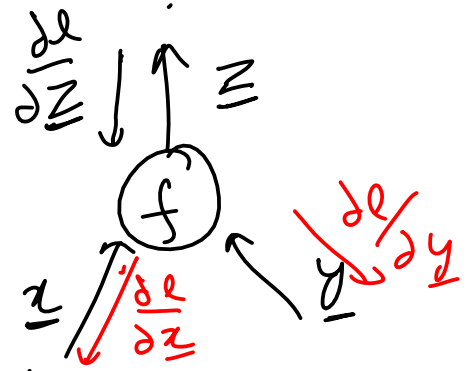graph is a scalar ( true for a loss
function )

$$\frac{\partial \ell}{\partial \underline{x}} = \frac{\partial \ell}{\partial \underline{z}} \frac{\partial f}{\partial \underline{x}} \leftarrow \text{Jacobian}$$

The responsibility of VJP function is to compute

$$\frac{\partial \ell}{\partial \underline{x}}, \frac{\partial \ell}{\partial \underline{y}} \text{ given } \frac{\partial \ell}{\partial \underline{z}} \leftarrow \text{vector}$$

and you can use $\underline{x}, \underline{y}$
$$\underline{f}(\underline{x}, \underline{y})$$

$\ell \in \mathbb{R}$

$$\frac{\partial \ell}{\partial \underline{z}} \downarrow \uparrow \underline{z}$$

$$\underline{x} \quad \frac{\partial \ell}{\partial \underline{x}} \qquad \frac{\partial \ell}{\partial \underline{y}} \quad \underline{y}$$

## Examples

$$\underline{f}(\underline{x}, \underline{y}) = \underset{n \times n}{I} \underline{x} + \underline{y}$$

Find the VJP

Assume an eventual scalar function

$$\frac{\partial}{\partial \underline{x}} \ell( \underline{f}(\underline{x}, \underline{y})) = \frac{\partial \ell}{\partial \underline{f}} \frac{\partial f}{\partial \underline{x}} = \frac{\partial \ell}{\partial \underline{f}} \left( I_{n \times n} + O_{n \times n} \right)$$

$$= \frac{\partial \ell}{\partial \underline{f}} = \frac{\partial \ell}{\partial \underline{z}}$$

$$\frac{\partial}{\partial \underline{y}} \ell(\underline{f}(\underline{x}+\underline{y})) = \frac{\partial \ell}{\partial \underline{f}} = \frac{\partial \ell}{\partial \underline{z}}$$

$$\frac{\partial A\underline{x}}{\partial \underline{x}} = A$$

$$\frac{\partial}{\partial \underline{x}} \begin{bmatrix} a_1^T \\ \vdots \\ a_n^T \end{bmatrix} \underline{x} = \begin{bmatrix} a_1^T \\ \vdots \\ a_n^T \end{bmatrix}$$

$$f(\underline{x}, \underline{y}) = \underline{x}^T \underline{y} \in \mathbb{R} \qquad \ell \in \mathbb{R}$$

$$\frac{\partial}{\partial \underline{x}} \ell(f(\underline{x}, \underline{y})) = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial \underline{x}} = \frac{\partial \ell}{\partial f} \underline{\underline{y}^T}$$

$$\frac{\partial}{\partial \underline{y}} \ell(f(\underline{x}, \underline{y})) = \frac{\partial \ell}{\partial f} \cdot \frac{\partial f}{\partial \underline{y}} = \frac{\partial \ell}{\partial f} \underline{x}^T$$

$$f(x) = \sin(x)$$

VJP?

$$f(x) = \frac{1}{1 + \exp(-x)}$$

VJP?

Project Convolution  x

$$f(W, \underline{I}) = W \otimes \underline{I}$$

VJP?

$$f(\underline{x}, \underline{w}) = \frac{1}{1 + \exp(\underline{w}^T \underline{x})}$$

VJP?