

## CS136 Lecture 2

### Outline

1. Random Variable & Probability
2. Joint, Conditional, Marginal
3. Sum Rule, Product Rule, Bayes Rule
4. Independence
5. Expectation, Mean and Variance

### Independence (conditional independence)

Intuitively,  $X$  is independent of  $Y$  if the probability of its value never depends on  $Y$ 's value.

$$\begin{aligned} p(X=x, Y=y) &= p(X=x | Y=y) \cdot p(Y=y) \\ &= p(X=x) \cdot p(Y=y) \end{aligned}$$

### Conditional Independence

We say  $X$  and  $Y$  are conditionally independent given  $Z$  if and only if

$$p(Y=y | X=x, Z=z) = p(Y=y | Z=z)$$

Again, if we know  $X$  and  $Y$  are conditionally independent

$$p(X, Y, Z) = p(X|Z) \cdot p(Y|Z) \cdot p(Z)$$

Proof:

$$\begin{aligned} p(X, Y, Z) &= p(X|Y, Z) \cdot p(Y|Z) \cdot p(Z) \quad (\text{product rule}) \\ &= \underbrace{p(X|Z)}_{X, Y \text{ is conditional independent}} \cdot p(Y|Z) \cdot p(Z) \end{aligned}$$

sum rule / product rule / Bayes rule

离散型

$$\text{sum rule: } p(x) = \sum_y p(x,y) = \sum_y p(x|y) \cdot p(y)$$

$$\text{product rule: } p(x,y) = \underline{p(x|y) \cdot p(y)}$$

连续型:

$$\text{sum rule: } p(x) = \int_y p(x,y) dy = \int_z p(x|y) \cdot p(y) dy$$

$$\text{product rule: } p(x,y) = p(x|y) \cdot p(y)$$

Bayes:

$$p(A|B) = \frac{p(B|A) \cdot p(A)}{\int_A p(B|A) \cdot p(A) dA}$$

推导:

$$\begin{aligned} p(x,y) &= p(x|y) \cdot p(y) \\ &= p(y|x) \cdot p(x) \end{aligned}$$

$$\Rightarrow \underbrace{p(y|x)}_{p(x,y)} \cdot p(x) = p(x|y) \cdot p(y)$$

$$\Rightarrow p(y|x) = \frac{p(x|y) \cdot p(y)}{p(x)}.$$

## Expectation - Mean and Variance

$$\text{Mean: } E(X) = \sum_{x \in \Omega} p(x) \cdot x$$

$$\text{Variance: } \boxed{\begin{aligned} \text{Var}(x) &= \sum_{x \in \Omega} E[(x - E(x))^2] \\ &= \sum_{x \in \Omega} p(x) \cdot (x - \mu)^2 \end{aligned}}$$

$$= \sum_{x \in \Omega} p(x) (x^2 - 2\mu x + \mu^2)$$

$$= \sum_{x \in \Omega} p(x)x^2 - \sum_{x \in \Omega} p(x) \cdot 2\mu x + \sum_{x \in \Omega} p(x) \cdot \mu^2$$

$$= \sum_{x \in \Omega} p(x)x^2 - 2\mu \underbrace{\sum_{x \in \Omega} p(x) \cdot x}_{\downarrow \mu} + \mu^2 \sum_{x \in \Omega} p(x)$$

$$= \sum_{x \in \Omega} p(x)x^2 - 2\mu^2 + \mu^2$$

$$\boxed{\text{Var}(x) = E(x^2) - \mu^2}$$

$$\boxed{\text{Var}(x) = E(x^2) - E(x)^2}$$

## CS13b Lecture 3

### Outline

1. Bernoulli Distribution
2. A spectrum of models for many coin flips
3. ML estimation for iid Bernoulli
4. Pros and Cons of Maximum Likelihood
5. Continuous Random Variables. PDFs and CDFs

Bernoulli Distribution.

$$P(X=x) = \begin{cases} \mu & x=1 \\ 1-\mu & x=0 \end{cases}$$

$$P(X=x) = \mu^x (1-\mu)^{1-x}$$

$X$  is Bernoulli Variable  $X \sim \text{Bern}(\mu)$

Some Properties:

1. Expectation.

$$\begin{aligned} E(X) &= \sum_{x \in \{0,1\}} x \cdot P(x|\mu) \\ &= 0 \cdot \mu^0 (1-\mu)^1 + 1 \cdot \mu (1-\mu)^0 \\ &= \mu \quad E[X] = \mu \end{aligned}$$

2. Variance

$$\text{Var} = E(X^2) - E(X)^2$$

$$\begin{aligned} E(X^2) &= \sum_{x \in \{0,1\}} x^2 P(x|\mu) \\ &= 0 \cdot \mu^0 (1-\mu)^1 + 1 \cdot \mu (1-\mu)^0 \\ &= \mu \end{aligned}$$

$$\text{Var} = \mu - \mu^2$$

$$\text{Var}(X) = \mu(1-\mu)$$

A spectrum of models for many coin flips

We perform  $N$  coin tosses  $\Rightarrow X_1, X_2, \dots, X_n$

We can model the joint of all tosses in several ways.

	joint Model	# parameter
Most general	$p(X_1, X_2, \dots, X_N)$	$N$ coins $\Rightarrow 2^N - 1$ free params
each toss is independent	$p(X_1)p(X_2) \dots p(X_N)$	$N$ one param for each toss.
each toss independent and identically distributed	$p(X_1  \mu)p(X_2  \mu) \dots p(X_N  \mu)$	1

Takeaways: Assumptions can simplify models, but need to carefully decide when appropriate

MLE for iid Bernoulli:

Given  $N$  observations of coin flip outcomes  $X_1, X_2, \dots, X_N$  where  $X_i \in \{0, 1\}$ . Assume an iid (independent & identically distribution model)

$$\begin{aligned} p(X_1 = X_1, X_2 = X_2, \dots, X_N) &= \prod_{i=1}^N \text{BernPMF}(X_i | \mu) \\ &= \prod_{i=1}^N \mu^{X_i} (1-\mu)^{1-X_i} \end{aligned}$$

maximum likelihood  $\Rightarrow$  Find a value of  $\mu$  that makes our observations most likely under assume model

$$\hat{\mu} = \arg \max_{\mu} \prod_{i=1}^N \mu^{X_i} (1-\mu)^{1-X_i} \quad \mu \in [0, 1]$$

Problem: inaccuracies in real application

the problem is that the product is too small.

Solution: use log-likelihood

why: Maximizing log-likelihood find the same optimal parameter  $\hat{\mu}$  as in L because log is a monotonic increasing function.

solve log-like likelihood

(单向逆传播)

$$\hat{\mu} = \text{argmax}_{\mu} \ln \prod_{i=1}^n \mu^{x_i} (1-\mu)^{1-x_i}$$

$$= \sum_{i=1}^n [x_i \ln \mu + (1-x_i) \ln (1-\mu)]$$

$$\frac{\partial L}{\partial \mu} = \sum_{i=1}^n \left[ \frac{1}{\mu} \cdot x_i + \frac{-1}{1-\mu} \cdot (1-x_i) \right]$$

$$= \frac{1}{\mu} \sum_{i=1}^n x_i - \frac{n}{1-\mu} + \frac{1}{1-\mu} \sum_{i=1}^n x_i \stackrel{\text{set}}{=} 0$$

$$= \frac{1}{\mu(1-\mu)} \sum_{i=1}^n x_i = \frac{n}{1-\mu}$$

$$\sum_{i=1}^n x_i = n\mu$$

$$\Rightarrow \boxed{\hat{\mu} = \frac{\sum_{i=1}^n x_i}{n}}$$

# of heads  
# total

Note.

for a constraint optimization problem, two ways to satisfy constraint

- (1) solve as if no constraint, see if answer satisfies
- (2) method of Lagrange multipliers

## Pros and Cons of Maximum Likelihood

### Advantages:

1. MLE are consistent

This means

- (1) our assumed likelihood distribution matches the true data-generating process
- (2) we have enough data

Suppose we pick some  $m^{\text{true}}$ , and then draw  $N$  observations  $X_n \sim p(x|m^{\text{true}})$   
As  $N \rightarrow \infty$ , we can prove that  $\hat{m}^{\text{ML}}(X_1, \dots, X_N) \rightarrow m^{\text{true}}$

2. MLE are equivalent to different parameterizations

### Limitations:

1. Overfitting (bad for small dataset)
2. Difficult of obtaining a solution  
(not all model have closed-form formula)
3. Solution is not always unique

## Continuous Random Variables. PDFs and CDFs

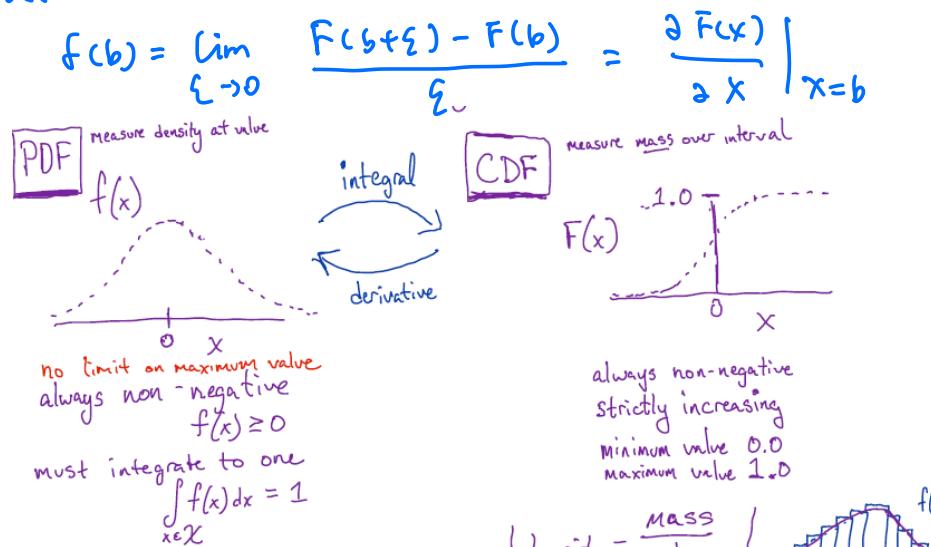
### Probability Density Functions

Let random variable  $X$  have continuous sample space  $\mathbb{R}$ . Define the cumulative distribution function  $F$  for  $X$  as

$$F(b) = P(X \leq b)$$

$$P(a \leq X \leq b) = F(b) - F(a)$$

We define the PDF of  $X$  as the function satisfying



Remember, PDFs indicate density. Think of integral as a Riemannian sum as # bins goes to infinity.

$$\int_{x \in X} f(x) dx \approx \lim_{\Delta x \rightarrow 0} \sum_{b=1}^B f(x_b) \Delta x = 1$$

density of bin  $b$       volume of bin  $b$       MASS

Density function can be larger than one. Implies very small volume. Example: Uniform over  $(0, \epsilon)$  has pdf  $\frac{1}{\epsilon}$

Let  $X \in \mathbb{R}$ ,  $Y \in \mathbb{R}$  be random variables

We can define joint PDF as  $P(X, Y)$ , which satisfies

- $P(X, Y) \geq 0$
- $\iint_{Y \in \mathbb{R}, X \in \mathbb{R}} P(X, Y) dxdy = 1$

Sum rule:  $P(X) = \int p(x, y) dy$  /  $P(Y) = \int p(x, y) dx$

Product rule:  $P(X, Y) = P(X|Y) \cdot P(Y)$   
 $= P(Y|X) \cdot P(X)$

## CS13b Lec4

Outline      Reading: 1, 2, 3, 2.1

1. Gamma functions
2. Beta Distributions
3. Beta-Bernoulli model and its posterior
4. MAP Estimation
5. Posterior Predictive Distributions

$$\begin{aligned}
 & \text{Prior / likelihood / posterior} \\
 & P(\mu) \quad P(x|\mu) \quad P(\mu|x) \\
 & P(x_1, \dots, x_N) = \int_0^1 P(\mu, x_1, \dots, x_N) d\mu \\
 & P(\mu|x_1, \dots, x_N) \quad (\text{posterior}) \\
 & = \frac{1}{P(x_1, \dots, x_N)} \cdot \underbrace{P(\mu)}_{\substack{\text{prior} \\ \downarrow}} \cdot P(x_1, \dots, x_N | \mu) \\
 & \qquad \qquad \qquad \uparrow \\
 & \text{Beta} \qquad \text{Bernoulli}
 \end{aligned}$$

Gamma Function  $P(x)$

- $(\alpha-1)! = P(\alpha) = \int_0^{+\infty} u^{\alpha-1} e^{-u} du$
- use  $\log P(x)$  in practice
- $P(x+1) = x \cdot P(x)$

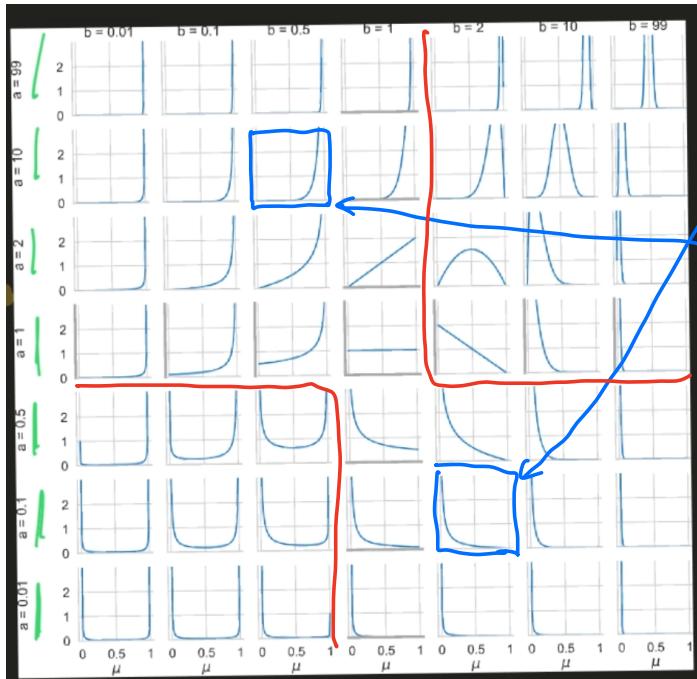
Beta Distributions

$$\begin{aligned}
 & \text{Beta PDF } P(\mu | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \quad 0 \leq \mu \leq 1 \\
 & \qquad \qquad \qquad \text{normalizing constant} \\
 & \qquad \qquad \qquad a > 0 \quad b > 0 \\
 & \qquad \qquad \qquad = C(a, b) \mu^{a-1} (1-\mu)^{b-1}
 \end{aligned}$$

- useful facts

$$\int_0^1 \text{Beta PDF}(\mu | a, b) d\mu = 1$$

$$\boxed{\int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \frac{1}{C(a, b)} = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}}$$



$a < 1, b < 1$ . U shape  
 $a \geq 1, b = 1$  — Uniform  
 $a < 1, b > 1$   
 $a > 1, b < 1$

Beta-Bernoulli model and its posterior

Define model over two random variables:

$$\mu \in [0, 1]$$

$$X_1, X_2, \dots, X_n \text{ with } X_i \in \{0, 1\}$$

Joint distribution?

$$p(X_1, \dots, X_N, \mu) = \frac{\prod_{n=1}^N \text{Bern PMF}(X_n | \mu)}{\text{likelihood}} \cdot \frac{\text{Beta PDF}(\mu | a, b)}{\text{"prior"} \quad p(\mu | X)} = \frac{p(X | \mu) \cdot p(\mu)}{p(X)}$$

Given joint model, we can ask about probabilities implied by model:

$$\text{Posterior } p(\mu | X_1, \dots, X_N)$$

after seeing  $N$  flip outcomes,  
what is likely value of  $\mu$   
what is the P of seeing  $N$   
flip outcome?

$$\text{Evidence } p(X_1, \dots, X_N)$$

$$\left. \begin{array}{l} \text{Bayes} \\ p(x) = \int p(x|\mu) d\mu \\ \text{sum rule} \end{array} \right\}$$

$$\text{Predictive Posterior } p(X_N | X_1, \dots, X_{N-1})$$

after seeing  $N-1$  flips, predict  
next outcome

$$\left. \begin{array}{l} \text{Bayes \& sum} \\ \text{rule} \end{array} \right\}$$

Posterior of  $\mu$  for Beta-Bernoulli

PMF  $\rightarrow$  discrete  $p(x_n | x)$   
 PDF  $\rightarrow$  continuous  $p(a \leq x \leq b)$

$$\begin{aligned}
 p(\mu | x_1, \dots, x_N) &= \frac{1}{p(x_1, x_2, \dots, x_N)} \frac{\prod_{n=1}^N \text{Bern PMF}(x_n | \mu)}{p(x|\mu)} \cdot \frac{\text{Beta PDF } (\mu | a, b)}{p(x)} \\
 &= \frac{1}{p(x_1, x_2, \dots, x_N)} \prod_{n=1}^N \mu^{x_n} ((1-\mu)^{1-x_n}) \cdot C(a, b) \mu^{a-1} (1-\mu)^{b-1} \\
 &= \frac{C(a, b)}{p(x_1, x_2, \dots, x_N)} \prod_{n=1}^N \mu^{\sum x_n} (1-\mu)^{\sum (1-x_n)} \cdot \mu^{a-1} (1-\mu)^{b-1} \\
 &= \text{Const} \cdot \frac{\mu^{\sum x_n + a-1}}{\mu^{\hat{a}-1} (1-\mu)^{\hat{b}-1}}, \text{ where } \hat{a} = a + T(x) \\
 &\quad \hat{b} = b + T(x)
 \end{aligned}$$

Posterior over  $\mu$  for Beta-Bern is Beta

A general pattern:

If you have an unknown distribution with pdf up to a constant

$$p(\theta | \lambda) = \text{const} \cdot f(\theta, \lambda)$$

and you know a distribution family  $D$  with pdf

$$p_D(\theta | \lambda) \approx l_0(\lambda) \cdot f(\theta, \lambda) \quad \text{if same as above}$$

$l_0$  is a known function

then your distribution must belong to family  $D$

MAP Estimation (Maximum a posteriori estimation)

meaning: most likely value of  $\mu$  by posterior density

Suppose we observe  $N$  flips outcomes,  $X_1, \dots, X_N$

(1) How to estimate a value for  $\mu$ ?

(2) How to predict the next flip  $X_{N+1}$ ?

Answer:

1) Find  $\hat{\mu}$  that is mostly likely under posterior  $p(\mu | x_1, x_2, \dots, x_n)$

2) Predict next flip using likelihood with  $\mu = \hat{\mu}$

$$p(X_{N+1}) = \text{Bern PMF}(X_{N+1} | \hat{\mu})$$

Optimization Problem: MAP

$$\begin{aligned}\hat{\mu} &= \arg \max_{\mu \in [0,1]} p(\mu | x_1, x_2, \dots, x_N) \\ &= \arg \max_{\mu \in [0,1]} \frac{1}{p(x_1, x_2, \dots, x_N)} \cdot p(\mu) \cdot p(x_1, x_2, \dots, x_N | \mu) \\ &= \frac{1}{p(x_1, x_2, \dots, x_N)} \log p(\mu) + \log p(x_1, x_2, \dots, x_N | \mu)\end{aligned}$$

$$= \arg \max_{\mu \in [0,1]} \log p(\mu) + \sum_{n=1}^N \log \text{Beta-Bern PMF}(x_n | \mu)$$

$$= \arg \max_{\mu \in [0,1]} \log p(\mu) + \log \mu^{\sum x_n} (1-\mu)^{\sum (1-x_n)}$$

:

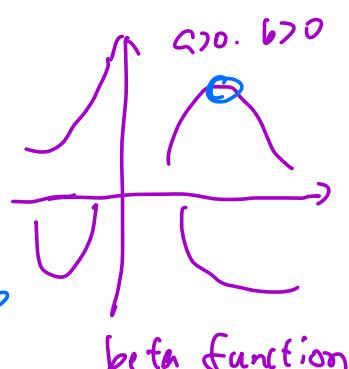
$$\Rightarrow \hat{\mu} = \frac{f(x) + a - 1}{f(x) + T(x) + a + b - 2}$$

$\underbrace{\phantom{f(x) + T(x) + a + b - 2}_{= N}}$

requires.

$$\begin{cases} f+a > 1 \\ T+b > 1 \end{cases}$$

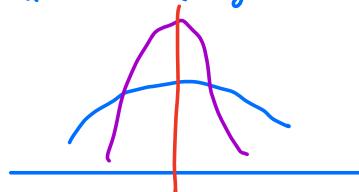
otherwise MAP  
doesn't exist



problems with MAP

(1) does not always exist

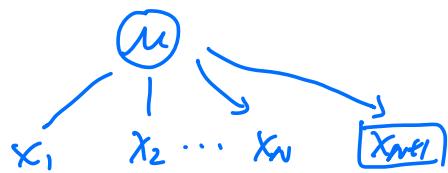
(2) Should we really condense whole contribution to one value?



these have same MAP!

But one is less confident

### Posterior Predictive Distributions



$$p(\mu) = \text{Beta}(\mu | a, b)$$

$$p(x_{N+1} | \mu) = \prod_n \text{Bern}(x_n | \mu)$$

$$p(x_{N+1} = 1 | x_1, \dots, x_N)$$

$$= \underbrace{\int p(x_{N+1} = 1, \mu | x_1, \dots, x_N) d\mu}_{\Psi}$$

$$p(x_{N+1}, \mu | x_1, \dots, x_N)$$

$$= p(x_{N+1} | \mu, x_1, \dots, x_N) \cdot \underbrace{p(\mu | x_1, \dots, x_N)}_{\substack{\uparrow \\ \text{Posterior given} \\ N \text{ coins}}}$$

$$= \underbrace{p(x_{N+1} | \mu)}_{\substack{\uparrow \\ \text{by iid} \\ p(x_{N+1} | \mu)}} \cdot \underbrace{\text{Bern}(x_{N+1} | \mu)}$$

$$= \int p(1 | \mu) \cdot \text{Beta}(\mu | \hat{a}, \hat{b}) d\mu$$

$$= C(\hat{a}, \hat{b}) \int \mu^{\hat{a}-1} (1-\mu)^{\hat{b}-1} d\mu$$

$$= C(\hat{a}, \hat{b}) \int \mu^{(\hat{a}+1)-1} (1-\mu)^{\hat{b}-1} d\mu$$

$$= C(\hat{a}, \hat{b}) \frac{1}{C(\hat{a}+1, \hat{b})}$$

$$= \frac{p(a+b)}{p(a)p(b)} \cdot \frac{p(a+1) \cdot p(b)}{p(a+1+b)}$$

1. Given  $p(a, b, c)$

$$p(a|c) = \int p(a, b|c) db$$

2. product rule

$$p(x_{N+1}, \mu | x_1, \dots, x_N) = \underbrace{p(x_{N+1} | \mu, x_1, \dots, x_N)}_{\substack{\text{joint} \\ \text{conditional}}} \cdot \underbrace{p(\mu | x_1, \dots, x_N)}_{\substack{\text{marginal}}}$$

$$p(x_T) = p(x) \cdot p(T)$$

$$p(x_n, \dots, x_1) = p(x_n | x_{n-1}, \dots, x_1) \cdot$$

$$\text{joint} = \text{marginal} \times \text{conditional} \quad p(x_{n-1}, \dots, 1)$$

$$\boxed{\int_0^1 \mu^{\hat{a}-1} (1-\mu)^{\hat{b}-1} d\mu = \frac{1}{C(\hat{a}, \hat{b})} = \frac{p(a)p(b)}{p(a+b)}}$$

$$\boxed{p(x+1) = x \cdot p(x)}$$

$$= \frac{p(a+b)}{p(a)} \cdot \frac{\hat{a} p(a)}{(\hat{a} + \hat{b}) - p(a+b)}$$

$$= \frac{\hat{a}}{(\hat{a} + \hat{b})}$$

Thus, by averaging over each possible  $\mu$ , weighted by its posterior density,

$$p(X_{n+1} = 1 | x_1, \dots, x_n) = \frac{\text{# head for } a}{N a + b}$$

Valid for  $a > 0$ .

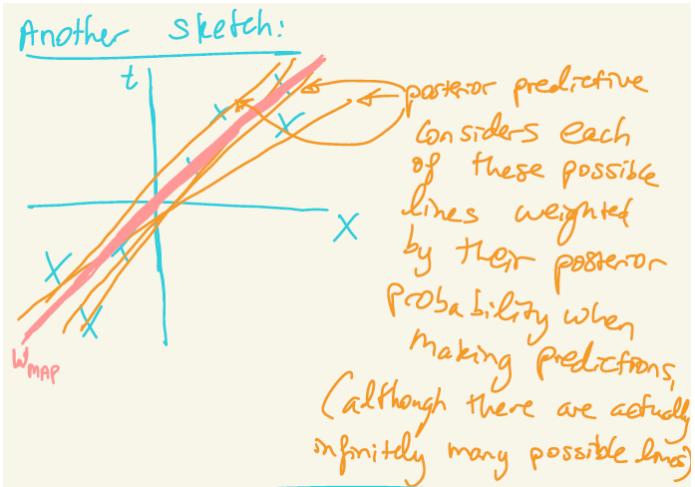
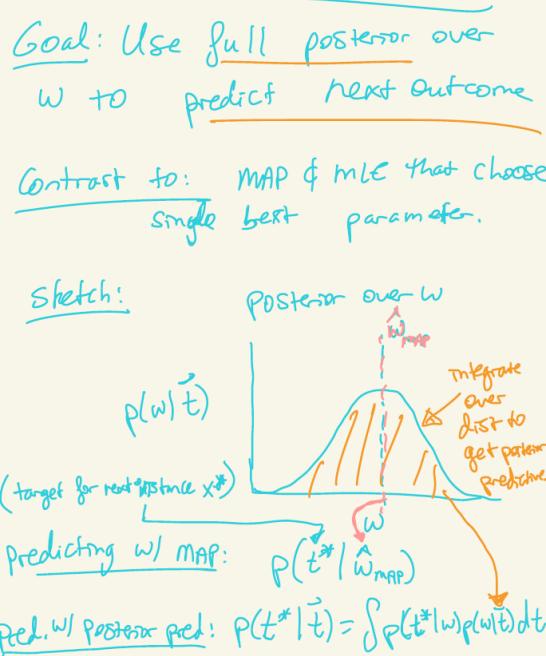
$b > 0$

$p(X_{n+1} | \hat{\mu})$  MAP

$$\int_0^1 p(X_{n+1} | \mu) p(\mu | x_1, \dots, x_n) d\mu \quad \text{PPE}$$

MLE	likelihood	Find max $\mu$
MAP	Posterior	Find max
Posterior prediction	Posterior	Integrate over $\mu$

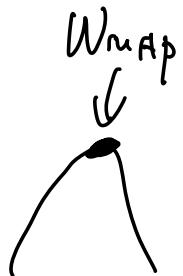
## 补充. MAP & Posterior for $\hat{y}^{*}$



$$p(w | t) = \frac{p(t(w) \cdot p(w))}{p(t)}$$

$$\propto p(t(w) \cdot p(w))$$

MAP:  $\arg \max_w \frac{p(t(w) \cdot p(w))}{\text{like}} \rightarrow \hat{w}_{MAP}$



$$\text{Probability} \propto p(t^* | X, t, \hat{w}_{MAP})$$



而 full bayesian

$$p(w | t) = p(t | w) \cdot p(w)$$

$$\text{Probability} \propto p(t^* | t)$$

$$\int p(t^* | w, t) \cdot p(w | t) dw = \int p(t^* | w, t) \cdot p(w | t) dw$$

$$= \underbrace{\int p(t^* | w) \cdot p(w | t) dw}_{\text{posterior}}$$

结论: Averaging should be better if we have limited data.

Bayesian Linear Regression  
shuhuai005 bilibili

Data:  $\{(x_i, y_i)\}_{i=1}^n, x_i \in \mathbb{R}^p, y_i \in \mathbb{R}$ .

Model:  $\begin{cases} f(x) = w^T x = x^T w \\ y = f(x) + \epsilon \\ \epsilon \sim N(0, \sigma^2) \end{cases}$

Bayesian Method:

参数:  $w$  不是随机变量  
 $w$  是一个概率分布

① Inference:  $p(w | \text{Data}) \rightarrow \text{posterior}$

$p(w | \text{Data}) \propto \text{likelihood} \times \text{prior}$

$\text{prior: } N(\mu_w, \Sigma_w)$   
 $\text{likelihood: } N(y | Xw, \sigma^2 I)$   
 $\text{posterior: } N(\mu_{\text{post}}, \Sigma_{\text{post}})$

② Prediction Given  $x^*, y^*$ :

$p(y^* | \text{Data}, x^*) = \int_w p(y^* | w, \text{Data}, x^*) p(w | \text{Data}, x^*) dw$

$p(y^* | w, x^*)$   
 $\downarrow$   
 $\text{posterior}$

# CS136 Lec 5 A

Outline Reading Sec 2.2

1. From Binary to Categorical Distributions
2. ML estimation for categorical Parameter
3. Dirichlet distribution
4. Dirichlet - Categorical model & its posterior
5. MAP Estimation for Dirichlet - Categorical

From Binary to Categorical Distributions

Binary Distribution  $\rightarrow$  2 outcomes

Categorical Distribution  $\rightarrow$   $V$  outcomes.

Binary (Bernoulli)

Sample space:  $\{0, 1\}$

Parameter:  $\mu \in [0, 1]$

PMF:  $\mu^x (1-\mu)^{1-x}$

Categorical Distribution

Sample space:  $\{1, 2, \dots, V\}$

Parameters:  $\mu = \{\mu_1, \dots, \mu_V\}$ ,

$$\begin{cases} \sum_{v=1}^V \mu_v = 1 \\ \mu_v \geq 0 \end{cases}$$

probability simplex

PMF:  $p(x=x|\mu) = \prod_{v=1}^V \mu_v^{x_v} = \mu_x$

notation:  $p(x=x|\mu) = \mu_x$

$x \in \{1, 2, \dots, V\}$

$$p(x=x|\mu) = \prod_{v=1}^V \mu_v^{x_v} = \mu_x$$

$\left\{ \begin{array}{ll} 1 & x=v \\ 0 & x \neq v \end{array} \right.$

one hot vector

Consider observing  $N$  words from vocabulary of size  $V$ .

Let random variables  $X_1, X_2, \dots, X_N$  indicate the words.

where  $X_n = [0 \ 0 \ 0 \ 0 \ | \ 0 \ 0 \ 0]$

$p(x|\mu) = \prod_{v=1}^V (\mu_v)^{x_v} (1-\mu_v)^{1-x_v}$

$\text{FREE: } p(x|\mu) = \prod_{v=1}^V \mu_v^{x_v}$

$$X_{nv} = \begin{cases} 1 & \text{if } n \text{ word is type } v \\ 0 & \text{otherwise} \end{cases}$$

Assume  $N$  words drawn independent & identically distributed (i.i.d) from Categorical with parameter  $\mu \in \Delta^V$

$$\begin{aligned} p(X_1, \dots, X_N | \mu) &= \prod_{n=1}^N \text{CatPMF}(X_n | \mu) \\ &= \prod_{n=1}^N \left( \frac{\mu_v}{\sum_v \mu_v} \right)^{X_{nv}} \mu_v^{1-X_{nv}} \quad \begin{matrix} \prod_{n=1}^N (\mu_v^{X_{nv}} \cdot \mu_v^{1-X_{nv}}) \\ = \prod_{n=1}^N \mu_v^{X_{nv} + 1 - X_{nv}} \end{matrix} \\ &= \frac{V}{\prod_{v=1}^V} \mu_v \underbrace{\sum_{v=1}^V X_{nv}}_{\substack{\text{number of times vocab} \\ \text{symbol } v \text{ appears in our} \\ \text{dataset}}} \\ &= \frac{V}{\prod_{v=1}^V} \mu_v \frac{m_v(x)}{m_v(x)} \quad m_v(x) = \sum_n X_{nv} \end{aligned}$$

ML estimation for categorical Parameter

$$\hat{\mu} = \arg \max_{\mu \in \Delta^V} \log p(X_1, X_2, \dots, X_N | \mu)$$

$$= \arg \max_{\mu \in \Delta^V} \log \left( \prod_{v=1}^V \mu_v^{m_v(x)} \right)$$

$$= \boxed{\arg \max_{\mu \in \Delta^V} \sum_{v=1}^V m_v(x) \cdot \log \mu_v}$$

$$\arg \max_{\mu \in \Delta^V} \sum_{v=1}^V m_v(x) \log \mu_v$$

$$\text{s.t. } \begin{cases} 1 - \sum_{v=1}^V \mu_v = 0 \\ \mu_v \geq 0 \end{cases} \quad \begin{matrix} \leftarrow \text{address with Lagrange multiplier} \\ \leftarrow \text{address with ignore, then check} \end{matrix}$$

Lagrange multiplier method

$$\text{Step 1. } L(\mu, \lambda) = \sum_{v=1}^V m_v(x) \log M_v + \lambda \left( 1 - \sum_{v=1}^V \mu_v \right) \quad (1 \neq 0)$$

$$\text{Step 2. } \frac{\partial}{\partial \mu_1} = 0 \Rightarrow \frac{m_1}{M_1} - \lambda = 0 \quad (1)$$

$$\frac{\partial}{\partial \mu_2} = 0 \Rightarrow \frac{m_2}{M_2} - \lambda = 0 \quad (2)$$

⋮

$$\frac{\partial}{\partial \mu_V} = 0 \Rightarrow \frac{m_V}{M_V} - \lambda = 0 \quad (V)$$

$$\frac{\partial}{\partial \lambda} = 0 \Rightarrow \underbrace{1 - \sum_{v=1}^V \mu_v}_{(V+1)} = 0$$

Summary (1) ... (V)

$$\sum_{v=1}^V \frac{m_v}{M_v} = \lambda$$

$$\Rightarrow \frac{1}{\lambda} \sum_{v=1}^V m_v = \sum_{v=1}^V \mu_v$$

plug in (V+1)

$$1 - \lambda \sum_{v=1}^V \mu_v = 0 \Rightarrow \frac{1}{\lambda} \sum_{v=1}^V m_v = 1$$

$$\lambda = N = \sum_{v=1}^V m_v$$

Return back to (1) ~ (V)

$$M_1^* = \frac{m_1}{\lambda} = \frac{m_1}{N}$$

$$M_2^* = \frac{m_2}{\lambda}$$

⋮

$$\Rightarrow M^* = \left[ \frac{m_1}{N}, \frac{m_2}{N}, \dots, \frac{m_V}{N} \right]$$

At last. check

$$\boxed{M^* \geq 0}$$

Proof

$$L(\mu, \lambda) = \sum_{v=1}^V m_v(x) \log M_v + \lambda \left( 1 - \sum_{v=1}^V \mu_v \right)$$

$$\frac{\partial L(\mu, \lambda)}{\partial \mu_v} = \frac{m_v}{M_v} - \lambda \stackrel{\text{set}}{=} 0$$

$$\frac{m_v}{M_v} = \lambda$$

$$\Rightarrow M_v^* = \frac{m_v}{\lambda}$$

由題意得

$$1 - \sum_{v=1}^V \mu_v = 0$$

$$\Rightarrow 1 - \sum_{v=1}^V \frac{m_v}{\lambda} = 0$$

$$M_v^* = \frac{m_v}{\lambda}$$

$\boxed{N}$

N is constant

## Dirichlet distribution

Dirichlet: prior

Random Variable  $\mu = [\mu_1, \mu_2, \dots, \mu_V]$

Sample space  $\mu \in \Delta^V$

Parameter  $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_V] \quad \alpha_1, 2, \dots, V > 0$

PDF:

$$\text{DirPPF}(\mu | \alpha) = \frac{\frac{P(\alpha_1 + \alpha_2 + \dots + \alpha_V)}{\prod_{v=1}^V P(\alpha_v)} \prod_{v=1}^V \mu_v^{\alpha_v - 1}}{\text{constant } c(\alpha)}$$

Can recognize that

$$\int_{\mu} \frac{P(\alpha_1 + \alpha_2 + \dots + \alpha_V)}{\prod_{v=1}^V P(\alpha_v)} \prod_{v=1}^V \mu_v^{\alpha_v - 1} d\mu = 1$$

$$\Rightarrow \int_{\mu} \mu_v^{\alpha_v - 1} = \frac{1}{c(\alpha)} = \frac{\prod_{v=1}^V P(\alpha_v)}{P(\alpha_1 + \alpha_2 + \dots + \alpha_V)}$$

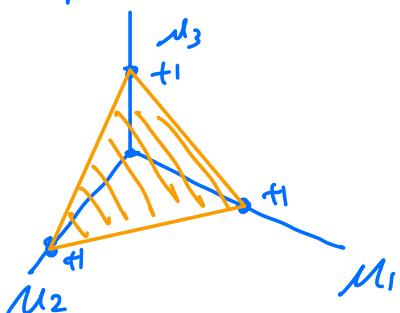
## Special cases & visualization

When  $V=2$ , Dirichlet reduces to Beta.

Beta

Sample space  $\mu \in [0,1]$

$$\text{PDF: } \frac{P(a+b)}{P(a)P(b)} \mu^{a-1} (1-\mu)^{b-1}$$

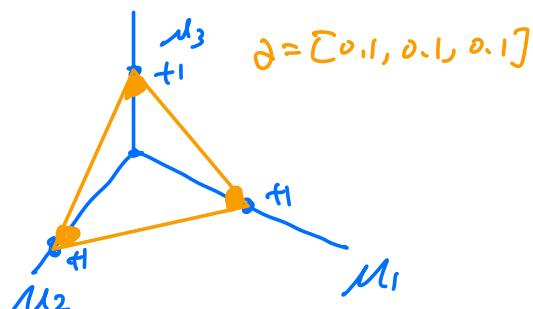


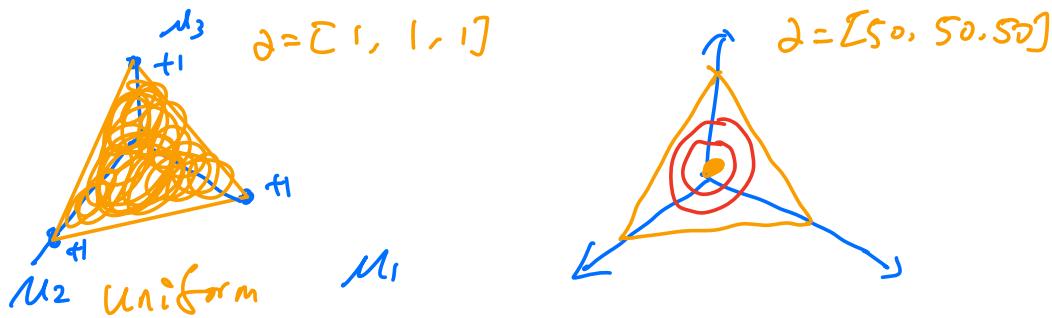
Dirichlet s.t.  $\mu_1 \geq 0$

$\mu_2 \geq 0$

$\mu_1 + \mu_2 = 1$

$$\frac{P(\alpha_1 + \alpha_2)}{P(\alpha_1)P(\alpha_2)} \cdot \mu_1^{\alpha_1 - 1} \cdot \mu_2^{\alpha_2 - 1}$$





Dirichlet - Categorical model & its posterior

Jointly Explain :

$[\mu_1, \mu_2, \dots, \mu_v]$  unknown probability vector  $\mu \sim \text{Dir}(\alpha)$

$x_1, x_2, \dots, x_N$   $N$  observed words  $x_n \stackrel{\text{iid}}{\sim} \text{Cat}(\mu)$

$$p(x_1, x_2, \dots, x_N, \mu) = \prod_{n=1}^N \text{Cat PMF}(x_n | \mu) \cdot \text{Dir PDF}(\mu | \alpha)$$

P.S. Beta - Bernoulli model

$$p(x_1, \dots, x_N, \mu) = \frac{\prod_{n=1}^N \text{Bern PMF}(x_n | \mu)}{\text{likelihood}} \cdot \frac{\text{Beta PDF}(\mu | a, b)}{\text{"prior"}}$$

Posterior of  $\mu$  under Dir-Cat

$$p(\mu | x_1, x_2, \dots, x_N)$$

$$= \frac{1}{p(x_1, x_2, \dots, x_N)} \cdot p(\mu) \cdot p(x_1, x_2, \dots, x_N | \mu)$$

$$= C_1 \cdot \text{Dir PDF}(\mu | \alpha) \cdot \prod_{n=1}^N \text{Cat PMF}(x_n | \mu) \quad \text{one-hot}$$

$$= C_1 \cdot \frac{P(\alpha_1 + \alpha_2 + \dots + \alpha_v)}{\prod_{v=1}^V P(\alpha_v)} \prod_{v=1}^V \mu_v^{\alpha_{v-1}} \cdot \prod_{n=1}^N \prod_{v=1}^V \mu_v^{x_{nv}}$$

$$= C_2 \cdot \prod_{v=1}^V \mu_v^{\alpha_{v-1}} \cdot \prod_{n=1}^N \prod_{v=1}^V \mu_v^{x_{nv}}$$

$$\begin{aligned}
 &= \prod_{V=1}^V M_V^{\alpha_{V-1}} \cdot \prod_{V=1}^V M_V \sum_{i=1}^N x_{iV} \\
 &= \prod_{V=1}^V M_V^{\alpha_{V-1}} \prod_{V=1}^V M_V^{m_V} \\
 &= \prod_{V=1}^V M_V^{\alpha_{V-1} + m_V} \quad \text{DirPDF}(\mu | \alpha) = \frac{P(\alpha_1 + \alpha_2 + \dots + \alpha_V)}{\prod_{V=1}^V P(\alpha_V)} \left[ \prod_{V=1}^V M_V^{\alpha_{V-1}} \right]
 \end{aligned}$$

Thus, posterior is Dirichlet Conjugate prior

$$PL(\mu | x_1, \dots, x_N) = \text{DirPDF}(\mu | \hat{\alpha}_1, \dots, \hat{\alpha}_V)$$

$$\text{where } \hat{\alpha}_v = \alpha_v + m_v$$

= prior pseudo count + count of symbol V

### MAP Estimation for Dirichlet - Categorical

$$\mu^{\text{MAP}} = \underset{\mu \in \Delta^V}{\arg \max} \log P(\mu | x_1, x_2, \dots, x_N)$$

$$= \underset{\mu \in \Delta^V}{\arg \max} \log \left( C(\hat{\alpha}) \prod_{V=1}^V M_V^{\hat{\alpha}_{V-1}} \right)$$

$$= \underset{\mu \in \Delta^V}{\arg \max} \sum_{V=1}^V (\hat{\alpha}_{V-1}) \log M_V$$

$$= \underset{\mu \in \Delta^V}{\arg \max} \sum_{V=1}^V (\hat{\alpha}_{V-1}) \log M_V$$

$$\mu_{\text{MAP}} = \left[ \frac{\hat{\alpha}_1-1}{S}, \frac{\hat{\alpha}_2-1}{S}, \dots, \frac{\hat{\alpha}_V-1}{S} \right]$$

$$S = \sum_v \hat{\alpha}_v - 1$$

## CS13b Lecture 5 B Lagrange multipliers

Outline

Appendix E

1. Recipe & example

2. Why it works?

Recipe & example

Given an objective function  $f$

$$[\mathbb{R}^D \rightarrow f \rightarrow \mathbb{R}]$$

$$\theta = [\theta_1, \theta_2, \dots, \theta_D]$$

Goal: Equality constrained optimal value

$$\theta^* = \underset{\theta}{\operatorname{argmin}} f(\theta)$$

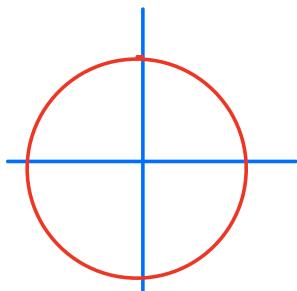
$$\boxed{s.t. g(\theta) = 0}$$

限制的方程

还有一个不等式  $g(\theta) \geq 0$ . 后面讲

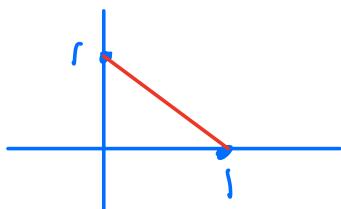
Example

Find  $\theta$  to maximize  $f$ , such that  $\theta$  lies on unit circle



$$g(\theta) = 1 - \theta_1^2 - \theta_2^2$$

Find  $\theta$  to maximize  $f$ , such that  $\theta$  sum to one.



$$g(\theta) = 1 - \sum_{d=1}^D \theta_d$$

Recipe

Step 1

Define objective  $f$  and constraint  $g$

Step 2

Define expanded objective, with new "multiplier"  $\lambda \neq 0$

$$L(\theta, \lambda) = f(\theta) + \lambda g(\theta)$$

Step 3  
calculus

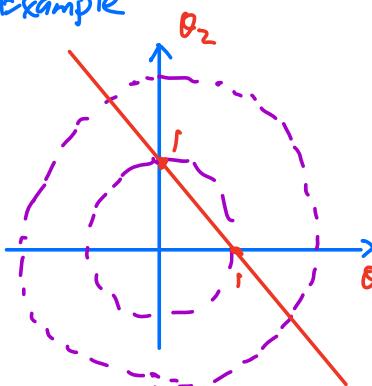
Setup system of  $D+1$  equations, where each partial derivative is zero

$$\frac{\partial}{\partial \theta_1} L = 0, \frac{\partial}{\partial \theta_2} L = 0, \dots, \frac{\partial}{\partial \theta_D} L = 0, \frac{\partial}{\partial \lambda} L = 0$$

Step 4  
algebra

Solve the system of equations to find optimal values of  $\theta_1^*, \theta_2^*, \dots, \theta_D^*, \lambda^*$

Example



step 1.

$$f(\theta) = 1 - \theta_1^2 - \theta_2^2 \leftarrow \text{objective}$$

$$g(x) = \theta_1 + \theta_2 - 1 \leftarrow \text{constraint}$$

step 2.

$$\begin{aligned} L(\theta, \lambda) &= f(x) + \lambda g(x) \\ &= 1 - \theta_1^2 - \theta_2^2 + \lambda(\theta_1 + \theta_2 - 1) \end{aligned}$$

step 3.

$$\frac{\partial L}{\partial \theta_1} = -2\theta_1 + \lambda \stackrel{\text{set}}{=} 0$$

$$\frac{\partial L}{\partial \theta_2} = -2\theta_2 + \lambda \stackrel{\text{set}}{=} 0$$

$$\frac{\partial L}{\partial \lambda} = \theta_1 + \theta_2 - 1 \stackrel{\text{set}}{=} 0$$

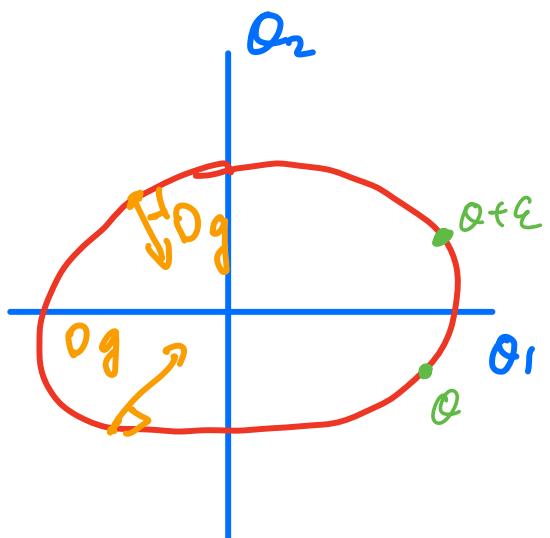
$$\Rightarrow \begin{cases} \theta_1 = \frac{1}{2} \\ \theta_2 = \frac{1}{2} \end{cases}$$

Why it works?

Consider the subset  $S$  of possible  $\theta$  values that satisfy the constraint.

Lemma 1: Any  $\theta$  that satisfies  $g(\theta)=0$  will have gradient perpendicular to constraint surface

$$\nabla g \perp S$$



Proof:

Examine two "close" points that satisfy  
call these points  $\theta$  and  $\theta + \varepsilon$

$$g(\theta) = 0 \quad (1)$$

$$g(\theta + \varepsilon) = 0 \quad (2)$$

Taylor:

$$g(\theta + \varepsilon) = g(\theta) + \varepsilon^T \nabla_{\theta} g(\theta)$$

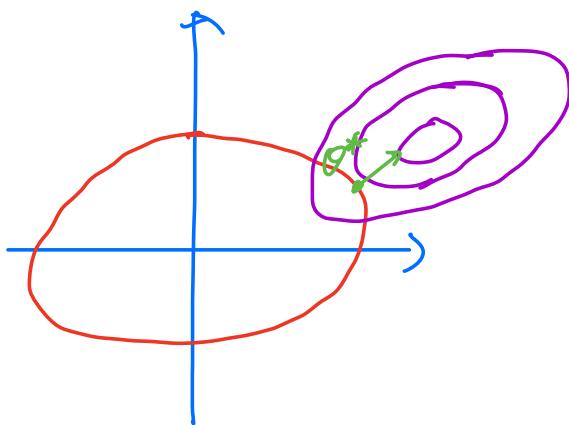
plug in (1):

$$g(\theta + \varepsilon) = \boxed{\varepsilon^T \nabla_{\theta} g(\theta) = 0}$$

Lemma 2:

At optimal  $\theta^*$ , the gradient of the objective  $f$  will be  
perpendicular to constraint surface

$$\nabla_{\theta} f|_{\theta=\theta^*} \perp \Sigma$$



punchline:

At optima  $\theta^*$ ,  $\nabla_{\theta} f$  and  $\nabla_{\theta} g$  are orthogonal to  $\Sigma$ .

thus they must be parallelled to each other

If  $\nabla_{\theta} f$ ,  $\nabla_{\theta} g$  are parallel, then  $\exists$  some  $\lambda \neq 0$  st.

$$\nabla_{\theta} f + \lambda \nabla_{\theta} g = 0$$

## Why it works (cont'd)? }

Thus, any optimal value  $\theta^*$  must satisfy these  $D+1$  equations

$\theta^*$  must be optimal

$$\vec{\nabla} f + \lambda \vec{\nabla} g = \vec{0}$$



$$\left[ \begin{array}{c} \frac{\partial}{\partial \theta_1} f \\ \vdots \\ \frac{\partial}{\partial \theta_D} f \end{array} \right] + \lambda \left[ \begin{array}{c} \frac{\partial}{\partial \theta_1} g \\ \vdots \\ \frac{\partial}{\partial \theta_D} g \end{array} \right] = \vec{0} \quad (1)$$

$$\left[ \begin{array}{c} \frac{\partial}{\partial \lambda} f \\ \vdots \\ \frac{\partial}{\partial \lambda} g \end{array} \right] = \vec{0} \quad (2)$$

$$\left[ \begin{array}{c} \frac{\partial}{\partial \theta_1} f \\ \vdots \\ \frac{\partial}{\partial \theta_D} f \end{array} \right] + \lambda \left[ \begin{array}{c} \frac{\partial}{\partial \theta_1} g \\ \vdots \\ \frac{\partial}{\partial \theta_D} g \end{array} \right] = \vec{0} \quad (D)$$

$\theta^*$  must satisfy constraints

$$g(\theta^*) = 0 \rightarrow \frac{\partial}{\partial \lambda} [\lambda g(\theta)] = 0 \quad (D+1)$$

Can view the whole system as partial derivatives  
of expanded objective

$$\alpha(\theta, \lambda) = f(\theta) + \lambda g(\theta)$$

where  $\theta \in \mathbb{R}^D$ ,  $\lambda \neq 0$  are unknown, to be solved for.