

Chain Rule

Bayes Rule

} read
these

Independence of RV

Expectation value of $f(x)$ w.r.t dist P

$$f: \Omega \rightarrow \mathbb{R}$$

Expectation value: most likely outcome of $f(x)$

$$E_p[f] = \sum_{x \in \Omega} p(x)f(x) \quad \text{discrete}$$

$$= \int_{\Omega} p f \quad \text{continuous}$$

E.g. $\Omega = \{\{i,j\} \mid 1 \leq i, j \leq 6\}$ two dice

$$f(\{i,j\}) = \begin{cases} i+j & i \neq j \\ 2i & i=j \end{cases}$$

$$E_p[f] = 7$$

Variance

$$\text{Var}_p(f(x)) = E_p \left[(f(x) - E_p(f(x)))^2 \right]$$

Fact: $\text{Var}_p(x) = E_p[x^2] - E_p[x]^2$

Proof: exercise

Standard Deviation:

$$SD_p(f(x)) = \sqrt{\text{Var}_p(f(x))}$$

(measure of clustering around mean)

Covariance:

$$\text{Cov}_p(f(x), g(x)) = E_p \left[(f(x) - E_p(f(x))) (g(x) - E_p(g(x))) \right]$$

measures of linear relationship

between f, g

(also depends on
scale of variables)

§ 3.9 Common Distributions

Bernoulli Distribution

$$\text{RV } X \quad \Omega = \{0, 1\} \quad \phi \in [0, 1]$$

$$P(X=0) = \phi \quad P(X=1) = 1-\phi$$

$$p(x) = \phi^x (1-\phi)^{1-x} \quad x=0,1$$

Exer: $E_p[x] = \phi$

$$\text{Var}_p(x) = \phi(1-\phi)$$

Multinoulli (Categorical)

$$\text{RV } X \quad \Omega = \{1, \dots, k\} \quad k \in \mathbb{N}$$

let $p_1, \dots, p_k \in [0,1]$ s.t. $\sum_i p_i = 1$

$$p(x=i) = p_i$$

Exer: Compute $E_p[x]$ and $\text{Var}(x)$

Gaussian (Normal) Distribution

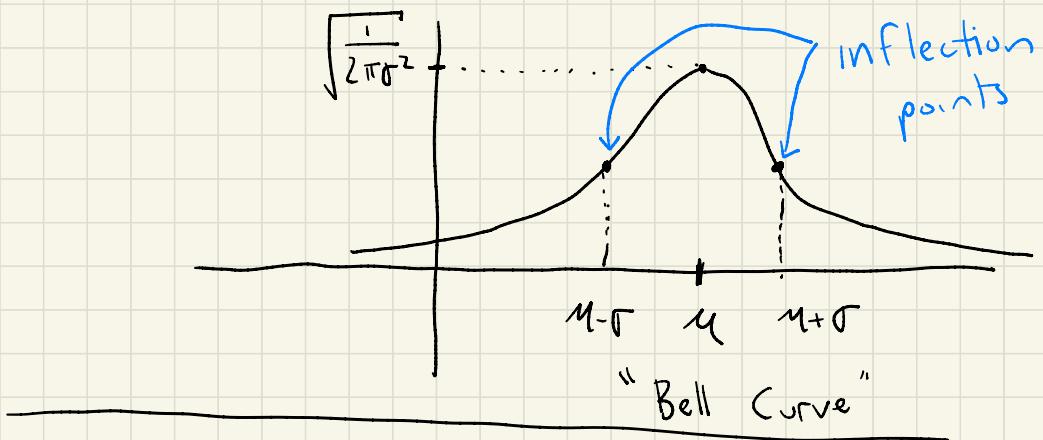
$$\text{RV } X \quad \Omega = \mathbb{R} \quad \text{fix } \mu \in \mathbb{R} \quad \sigma \in \mathbb{R}_{>0}$$

$$p(x=x) = N(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

Exer: $E_p[x] = \mu$

$$\text{Var}_p(x) = \sigma^2$$

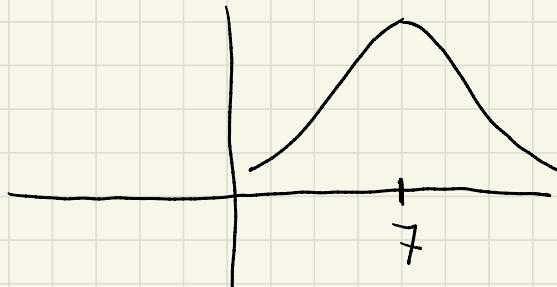
$$\int_{\mathbb{R}} N(x; \mu, \sigma^2) dx = 1$$



Def $\exp(x) = e^x$

Central Limit Theorem Sum of many independent random variables is approximately normally distributed.

E.g. 2-dice dist of $i+j$



Multivariate Gaussian

$$\Omega = \mathbb{R}^n \quad \mu \in \mathbb{R}^n \quad \Sigma \in \mathbb{R}^{n \times n} \quad \begin{matrix} \text{pos def} \\ \text{symmetric} \end{matrix}$$

$$N(\vec{x}; \mu, \Sigma) = \frac{1}{(2\pi)^n \det(\Sigma)} \exp\left(-\frac{1}{2} (\vec{x} - \mu)^T \Sigma^{-1} (\vec{x} - \mu)\right)$$

Σ covariance matrix (often
diagonal)

Others: exponential, Laplace, dirac, empirical
(read in book)

Mixing Distributions

Given distributions labeled by $c=1, \dots, k$

and a categorical dist on $1, \dots, k$

$$\text{define } p(x) = \sum_i P(c=i) P(x|c=i)$$

\uparrow \uparrow
 Categorical distribution i

Common to mix Gaussian distributions

Chapter 4: Optimization / numerical computation

§4.1 Rounding error (overflow / underflow)

§4.2 Poor conditioning

$$\text{matrix } A \quad \text{cond}(A) = \max_{i,j} \left| \frac{\lambda_i}{\lambda_j} \right|$$

If this is too big, you may run into numerical problems when computing with A.

§4.3 Gradient Based Optimization

In machine learning, often have

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \quad (\text{or } f: U \rightarrow \mathbb{R}, U \subseteq \mathbb{R}^n)$$

and we want

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

f: objective / criterion / cost / loss / error function

Recall Given $\hat{u} \in \mathbb{R}^n$ unit vector ($\hat{u}^\top \hat{u} = 1$)

directional derivative

$$D_{\hat{u}} f(\vec{x}) = \lim_{h \rightarrow 0} \frac{F(\vec{x} + \hat{u} h) - f(\vec{x})}{h}$$

$$= \left. \frac{\partial}{\partial \alpha} F(\vec{x} + \alpha \hat{u}) \right|_{\alpha=0}$$

$$\frac{\partial}{\partial x_i} f(\vec{x}) = D_{e_i} f(\vec{x}) \quad e_i = (0 \dots \underset{i}{1} \dots 0)^T$$

$$\nabla f(\vec{x}) = \left(\frac{\partial F}{\partial x_1}(\vec{x}) \dots \frac{\partial F}{\partial x_n}(\vec{x}) \right)^T$$

Question Which direction \hat{u} maximizes

$$D_{\hat{u}} f(\vec{x}) ?$$

Answer: $D_{\hat{u}} f(\vec{x}) = \left. \frac{\partial}{\partial \alpha} F(\vec{x} + \alpha \hat{u}) \right|_{\alpha=0}$

$$= \hat{u}^\top \nabla_x f(\vec{x}) \quad (\text{chain rule})$$

$$= \| \nabla_x f(\vec{x}) \| \cos \theta$$

maximized when $\hat{u} = \frac{\nabla_x f(\vec{x})}{\|\nabla_x f(\vec{x})\|}$

So $\nabla_x f(\vec{x})$ is the direction of max inc
 $-\nabla_x f(\vec{x})$ " " " " " max dec

Gradient Descent Algorithm

To find $\vec{x}^* = \operatorname{argmin} f(\vec{x})$ (hopefully)

Pick \vec{x} randomly and choose parameters

$\alpha > 0$ (step size)

$\varepsilon > 0$ (stopping signal)

While True:

$$\text{prev} = \vec{x}$$

$$\vec{x} = \vec{x} - \alpha \nabla_x f(\vec{x})$$

IF $\|\vec{x} - \text{prev}\|_2 < \varepsilon$:

break

Note: May want to change value of α at each step, and play with value of ϵ .