

Conditional probability $P(A|B)=\frac{P(A \cap B)}{P(B)}$: Probability for variable, given other variable \circ If X, Y are discrete:
 $p(x|y)=\frac{p(x,y)}{p(y)}$ \circ If X, Y are continuous:
 $f(x|y)=\frac{f(x,y)}{f(y)}$ \circ If X is discrete, Y is continuous:
 $p(x|y)=\frac{f(x,y)}{f(y)}$ \circ If X is continuous, Y is discrete:
 $f(x|y)=\frac{f(x,y)}{p(y)}$ \circ Properties: $\blacksquare P(A|B)=1-P(A^c|B)$
 $\blacksquare P(A_1|B)+P(A_2|B)+...=1$ \blacksquare If conditioning on subset S :
 $p(x|y)=\frac{p(x,y)}{p(y)} \times \frac{p(y \in S)}{p(S)}$ \circ $x \in S$ Bayesian
 $p(x|S)=\begin{cases} 0 & x \notin S \\ \text{Bayesian} & x \in S \end{cases}$
terminology: \circ Prior P (parameter) \circ Posterior P (parameter|data) \circ Likelihood P (data|parameter) \circ Evidence P (data) \circ Bayes theorem:
Posterior $P(A|B)$ = Likelihood $P(B|A)$ \times Prior $P(A)$ where Evidence $P(B)$
 $P(B)$ can be rewritten in marginalized form over A
Measures n^{th} moment $= \mathbb{E}(X^n) = \int_{-\infty}^{\infty} x^n f(x) dx$
Expected value \rightarrow Generally, \circ If X is discrete:
 $\mathbb{E}(X)=\sum_{x \in X} x \times p(x)$ \circ If X is continuous:
 $\mathbb{E}(X)=\int_{-\infty}^{\infty} x \times f(x) dx$ \circ If Y is discrete:
 $\mathbb{E}[X]=\sum_{y \in Y} \mathbb{E}[X|Y=y] p(y)$ \circ If Y is continuous:
 $\mathbb{E}[X]=\int_{-\infty}^{\infty} \mathbb{E}[X|Y=y] f(y) dy$ For functions: \circ $g(X)$ is a function, \circ If X is discrete: $\mathbb{E}(g(X))=\sum_{x \in X} g(x) \times p(x)$ \circ If X is continuous: $\mathbb{E}(g(X))=\int_{-\infty}^{\infty} g(x) \times f(x) dx$ For probabilities:
 \circ Count as functions: A is an event, X is a random variable \circ If X is discrete: $\mathbb{E}[p(X|A)]=\sum_{x \in X} p(x|A) p(x)$ \circ If X is continuous: $\mathbb{E}[p(X|A)]=\int_{-\infty}^{\infty} p(x|A) f(x) dx$ \circ If X is discrete: $\mathbb{E}[p(X|A)]=\sum_{x \in X} p(x|A) \times \sum_{x \in X} p(A, x) = p(A)$ \circ If X is continuous: $\mathbb{E}[p(X|A)]=\int_{-\infty}^{\infty} p(A, x) f(x) dx = p(A)$ For conditions: $\bullet A$ is an event, X is a random variable \circ If X is discrete: $\mathbb{E}(X|A)=\sum_{x \in X} x \times p(x|A)$ \circ If X is continuous: $\mathbb{E}(X|A)=\int_{-\infty}^{\infty} x \times f(x|A) dx$ \circ $\mathbb{E}(A|X)=P(A|X)$ For vectors: \circ Expectation of a vector is the expectation of each of its elements \circ If X is discrete:
 $\mathbb{E}(x)=\sum_{x_1} \dots \sum_{x_n} x_i p(x_1, \dots, x_n) = \mu$ \circ If X is continuous: $\mathbb{E}(x)=\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} x^T f_{X_1, \dots, X_n}(x_1, \dots, x_n) dx_1 \dots dx_n = \mu$ Properties: \circ $\mathbb{E}(a) = a$ \circ $\mathbb{E}(aX+b) = a\mathbb{E}(X) + b$ \circ $\mathbb{E}(aX+bY) = a\mathbb{E}(X) + b\mathbb{E}(Y)$ \bullet For independent variables: \circ $\mathbb{E}(XY) = 0$ \circ $\mathbb{E}(X+Y)^2 = \mathbb{E}(X^2) + \mathbb{E}(Y^2)$ \bullet For independent variables: \circ $\mathbb{E}(X|Y) = \mathbb{E}(X)$ \circ $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$ \bullet For vectors: If $y = Ax$: $\mathbb{E}(Ax) = A\mathbb{E}(x)$ \circ $\mathbb{E}[B(X|Y)] = \mathbb{E}(X)$ \bullet Cauchy Schwarz ineq.: $\mathbb{E}(X, Y)^2 \leq \mathbb{E}(X^2)\mathbb{E}(Y^2)$ Median \rightarrow Real number M defined by $P(X < M) = P(X > M)$ Standard deviation $= \sqrt{V(X)}$ Covariance \rightarrow Univariate variance of a random variable: $V(X) = \mathbb{E}(X - \mathbb{E}(X))^2 = \mathbb{E}(X^2) - \mathbb{E}(X)^2$ where $\mathbb{E}(X^2)$ is the unnormalized correlation resp. inner product \bullet Univariate covariance of two random variables: $Cov(X, Y) = \mathbb{E}((X - \mathbb{E}(X))(Y - \mathbb{E}(Y))) = \mathbb{E}(XY) - \mu_X \mu_Y$ where $\mathbb{E}(XY)$ is the unnormalized correlation resp. inner product \bullet Proof (heuristically for variance): $V(X) = \mathbb{E}(X - \mathbb{E}(X))^2 = \mathbb{E}[X^2 - 2X\mathbb{E}(X) - X\mathbb{E}(X) + \mathbb{E}(X)^2] = \mathbb{E}[X^2] - \mathbb{E}[X]\mathbb{E}(X) - \mathbb{E}[X]\mathbb{E}(X) + \mathbb{E}(X)^2 = \mathbb{E}(X^2) - \mathbb{E}(X)^2$ is the second moment \bullet Multivariate covariance matrix of a vector: $\Sigma = Cov(X) = \mathbb{E}((X - \mathbb{E}(X))(X - \mathbb{E}(X))^T) = \mathbb{E}(xx^T) - \mathbb{E}(x) \mathbb{E}(x)^T = R - \mu_X \mu_X^T$ $\left[\begin{array}{c} \mathbb{E}(x_1^2) \mathbb{E}(x_1 x_2) \\ \mathbb{E}(x_2 x_1) \mathbb{E}(x_2^2) \end{array} \right] - \left[\begin{array}{c} \mathbb{E}(x_1)^2 \mathbb{E}(x_1) \mathbb{E}(x_2) \\ \mathbb{E}(x_2) \mathbb{E}(x_1) \mathbb{E}(x_2) \mathbb{E}(x_2^2) \end{array} \right] = \left[\begin{array}{c} V(x_1) Cov(x_1, x_2) \\ Cov(x_2, x_1) V(x_2) \end{array} \right]$ where $R = \mathbb{E}(xx^T)$ is the unnormalized correlation matrix $\bullet \Sigma$ and R are symmetric and psd Properties - variance: $V(a) = a^T V(X) a$ $\bullet V(aX+b) = a^T V(X) a$ $\bullet V(X+Y) = V(X) + 2Cov(X, Y) + V(Y)$ \bullet For uncorrelated (and independent) variables: $V(X+Y) = V(X) + V(Y)$ \bullet For independent variables: $V(X|Y) = \mathbb{E}((X - \mathbb{E}(X))^2 | Y) = V(X)^2$ \bullet For vector $y = Ax$: $V_y = A V_X A^T$ \bullet For zero-mean variable: $V(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2 = \mathbb{E}(X^2)$ since $\mathbb{E}(X) = 0$ $\bullet V[X] = V[\mathbb{E}[X|Y]] + \mathbb{E}[V[X|Y]]$ $\bullet V[\sum_i x_i] = \sum_i V[x_i] + 2 \sum_{i < j} Cov[x_i, x_j]$ Properties - covariance: $Cov(X, X) = V(X)$ $\bullet Cov(aX+bY, Z) = aCov(X, Z) + bCov(Y, Z)$ \bullet If covariance of 2 random variables is 0 resp. $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$, they are uncorrelated, but not necessarily independent \bullet If unnormalized correlation of 2 random variables is 0 resp. $\mathbb{E}(XY) = 0$, they are orthogonal, but not necessarily independent \bullet For vector $y = Ax$: $\Sigma_y = A \Sigma_X A^T$ $\bullet R_y = A R_X A^T$ \bullet For zero-mean variables: $Cov(X, Y) = \mathbb{E}(XY) - \mu_X \mu_Y = \mathbb{E}(X, Y)$ since $\mu_X = \mu_Y = 0$ \bullet Cauchy Schwarz ineq.: $Cov(X, Y)^2 \leq V(X)V(Y)$ Correlation \rightarrow Normalized covariance \bullet Univariate correlation of a random variable: $Cor(X, Y) = \frac{Cov(X, Y)}{\sqrt{V(X)}\sqrt{V(Y)}}$ \bullet Multivariate correlation matrix of a vector: $\rho = Cor(X)$ $\bullet \rho$ is symmetric and psd \bullet Correlation is bounded between 0 and 1, given Cauchy Schwarzby: $\rho \rightarrow 0$, then $(\hat{S} - \mathbb{E}X \hat{S} \geq \epsilon) \rightarrow 0$ \bullet Absolute deviation given by $\rho(P(|\hat{S}_n - \mathbb{E}X \hat{S}_n| \geq \epsilon)) \leq \rho(P(\hat{S}_n - \mathbb{E}X \hat{S}_n \geq \epsilon))$

Inequality \circ If correlation of two random variables is 0, they are not necessarily independent
Probability Distributions Normal distribution $\rightarrow X \sim N(\mu, \sigma^2)$
For univariate, PDF: $\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2} (x-\mu)^2}$
For multivariate, PDF: $\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$
Convolution: $f \otimes g(x) = \int_{-\infty}^{\infty} f(x-y)g(y)dy$
Standard normal distribution \rightarrow Normal distribution, standardized with z-score $z = \frac{x-\mu}{\sigma}$, which results in $\mu=0$ and $\sigma=1$
Bernoulli distribution \rightarrow Bernoulli distribution (probability p) or failure (probability $1-p$) $\bullet X \sim \text{Bernoulli}(p)$ \bullet PDF: $p(x) = p^x (1-p)^{1-x}$ \bullet Mean: $\mathbb{E}(x) = p$ \bullet Variance: $V(x) = p(1-p)$ \bullet Binomial distribution $\rightarrow n$ independent Bernoulli trials with k successes $\bullet X \sim \text{Bin}(n, p)$ \bullet PDF: $\binom{n}{k} p^k (1-p)^{n-k}$ \bullet Mean: $\mathbb{E}(x) = np$ \bullet Variance: $V(x) = np(1-p)$ \bullet Poisson distribution \rightarrow $\bullet X \sim \text{Pois}(\lambda)$ \bullet PDF: $e^{-\lambda} \frac{\lambda^x}{x!}$ \bullet Mean: $\mathbb{E}(x) = \lambda$ \bullet Variance: $V(x) = \lambda$ \bullet Beta distribution $\rightarrow X$ takes values in $[0, 1]$ \bullet Represents the probability of a Bernoulli process after observing α successes and $\beta-1$ failures $\bullet X \sim \text{Beta}(\alpha, \beta)$ where $\alpha, \beta > 0$ \bullet PDF: $\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$ where $\Gamma(\alpha) = \int_0^\infty u^{\alpha-1} e^{-u} du$ \bullet Mean: $\mathbb{E}(x) = \frac{\alpha}{\alpha+\beta}$ \bullet Variance: $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$ \bullet Dirichlet distribution $\rightarrow X$ takes values in $[0, 1]^n$ \bullet Multivariate extension of Beta distribution \bullet $Dir(x|\alpha) = \frac{1}{B(\alpha)} \prod_{i=1}^n x_i^{\alpha_i-1}$, where $B(\alpha)$ is the multivariate generalization of the Beta function: $B(\alpha) = \frac{\prod_{i=1}^n \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^n \alpha_i)}$ \bullet Uniform distribution \rightarrow Assume x is uniformly distributed between $[a, b]$ \bullet PDF: $f(x) = \frac{1}{b-a}$ if $a \leq x \leq b$, else 0 \bullet CDF: $F(x) = \frac{x-a}{b-a}$ if $a \leq x \leq b$, 1 if $x > b$, 0 if $x < a$
Other Concepts Law of large numbers \rightarrow Sample mean of iid variables converges to population mean as $n \rightarrow \infty$ \bullet Weak law of large numbers: $\lim_{n \rightarrow \infty} P(|\bar{M}_n - \frac{1}{n} \sum_{k=1}^n X_k| < \epsilon) = 1$ \bullet Strong law of large numbers: $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n X_k = \mu_X$ with probability 1
Union bound $\rightarrow P(\cup_i A_i) \leq \sum_i P(A_i)$
Jensen's ineq. \rightarrow Relates expected value of a convex func. of a random variable to the convex func. of the expected value of that random variable $\mathbb{E}(f(X)) \geq f(\mathbb{E}(X))$
Markov's ineq. $\rightarrow p(x \geq t) \leq \frac{\mathbb{E}(x)}{t}$. Interesting only for $t \geq \mathbb{E}(x)$ because $p(x \geq t)$ must then be less than or equal to 1.
Generalizations: $\bullet p(|x| \geq t) \leq \frac{\mathbb{E}(|x|)}{t}$
 $\bullet p(|x| \geq t) \leq \frac{\mathbb{E}(x^2)}{t^2}$
Hoeffding's Lemma \rightarrow For random variable with $\mathbb{E}[x]=0$, and $a \leq x \leq b$, and $s > 0$: $\mathbb{E}[e^{sX}] \leq \exp(\frac{s^2(b-a)^2}{8})$
Hoeffding's inequality \rightarrow For random variables x_i that fall in the interval $[a_i, b_i]$ with probability 1, and $S_n = \sum_{i=1}^n x_i$, and $t > 0$: $P(S_n - \mathbb{E}X S_n \geq t) \leq \exp(-\frac{t^2}{\sum_{i=1}^n (b_i - a_i)^2})$
 $\bullet P(S_n - \mathbb{E}X S_n \leq -t) \leq \exp(-\frac{t^2}{\sum_{i=1}^n (b_i - a_i)^2})$ Proof:
 \bullet Consider the probability $P(S_n - \mathbb{E}[S_n] \geq t)$ \bullet Using Markov's ineq.: $P(S_n - \mathbb{E}[S_n] \geq t) = P(e^{s(S_n - \mathbb{E}[S_n])} \geq e^{st}) \leq \frac{\mathbb{E}[e^{s(S_n - \mathbb{E}[S_n])}]}{e^{st}}$ \bullet Using independence of X_1, \dots, X_n : $\mathbb{E}[e^{s(S_n - \mathbb{E}[S_n])}] = \prod_{i=1}^n \mathbb{E}[e^{s(X_i - \mathbb{E}[X_i])}]$ \bullet For each term, we use the fact that $X_i \in [a_i, b_i]$, and apply the lemma ineq.: $\prod_{i=1}^n \mathbb{E}[e^{s(X_i - \mathbb{E}[X_i])}] \leq \prod_{i=1}^n \exp(\frac{s^2(b_i - a_i)^2}{8})$ \bullet Plugging this back in: $e^{-st} \times \prod_{i=1}^n \exp(\frac{s^2(b_i - a_i)^2}{8}) = e^{-st} \times \exp(\frac{s^2}{8} \sum_{i=1}^n (b_i - a_i)^2) = P(S_n - \mathbb{E}[S_n] \geq t) \leq \exp(-\frac{st}{\sum_{i=1}^n (b_i - a_i)^2})$ \bullet If we set $s = \frac{4t}{\sum_{i=1}^n (b_i - a_i)^2}$ to min. the bound, we get: $P(S_n - \mathbb{E}[S_n] \geq t) \leq \exp(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2})$ \bullet Similarly for $P(S_n - \mathbb{E}[S_n] \leq -t)$. In the special case of normalized sums of iid variables, where $\hat{S}_n = S_n/n$ and $t = n\epsilon$: Δ Delta given by: $P(\hat{S}_n - \mathbb{E}X \hat{S}_n \geq \epsilon) \leq \exp(-\frac{2n\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2/n})$ \bullet As $n \rightarrow \infty$, then $(\hat{S}_n - \mathbb{E}X \hat{S}_n \geq \epsilon) \rightarrow 0$ \bullet Absolute deviation given by $\rho(P(|\hat{S}_n - \mathbb{E}X \hat{S}_n| \geq \epsilon)) \leq \rho(P(\hat{S}_n - \mathbb{E}X \hat{S}_n \geq \epsilon))$

$\epsilon \sqrt{S_n} - \mathbb{E}X \hat{S}_n \leq \epsilon$ \bullet By the union bound: $P(\hat{S}_n - \mathbb{E}X \hat{S}_n \geq \epsilon) + P(\hat{S}_n - \mathbb{E}X \hat{S}_n \leq -\epsilon)$ \bullet By Hoeffding's ineq.: $\leq 2 \exp(-\frac{2n\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2/n})$
Chebyshev's ineq. $\rightarrow p(|x - \mu_X| \geq \alpha |\sigma_X|) \leq \frac{1}{\alpha^2}$ resp. $p(|x - \mu_X| \geq \alpha) \leq \frac{1}{\alpha^2}$. Interesting only for $\alpha > 1$. Implications:
 \bullet For n variables: $p(|(S_n - \mu_X)/n| \geq \epsilon) \leq \frac{\sigma_X^2}{n\epsilon^2}$ where $S_n = \frac{1}{n} \sum_{i=1}^n X_k$ is the sample mean
Sufficient statistics $\rightarrow Z = g(X)$ is a sufficient statistic for estimating X if X can be estimated as well from Z as from Y , i.e. condensing Y to Z does not entail any loss of information about X \bullet Conditioned on Z , Y is independent of X : $p(Y|Z, X) = p(Y|Z)$ \bullet For sufficient statistics, the MLE of X from Y is the same as the MLE of X from Z : $p_{\text{argmax}} p(Y|X)p(Y) = \text{argmax}_X p(Z|X)p(Y)$ $\bullet p(X|Z) \propto p(X|Y)$
Hypothesis Testing Terminology \rightarrow \bullet Hypothesis: $\circ H_0$: Accepted null hypothesis, e.g. $p = p_0$, $p_1 - p_2 = p_0$, $1 - p_0, 2 = 0$ $\circ H_A$: Alternative hypothesis, e.g. $p \neq p_0$, $p_1 - p_2 \neq p_0$, $1 - p_0, 2 \neq 0$ \bullet Test types: \circ Two-sided: $p_1 - p_2 \neq p_0$, $H_A: p \neq p_0$ \circ One-sided upper tail: $H_0: p \geq p_0$, $H_A: p > p_0$ \circ One-sided lower tail: $H_0: p \leq p_0$, $H_A: p < p_0$ \bullet Errors: \circ True positive: Chose H_0 , and H_0 obtains \circ False positive, type I error: Chose H_0 , but H_0 obtains \circ True negative: Chose H_A , and H_A obtains \circ False negative, type II error: Chose H_0 , but H_A obtains \bullet Significance level α : $\circ \alpha \geq p$ (type I error) $= p(\bar{x} \geq z | H_0) = p(z_n \geq z_\alpha | H_0)$ with eq. for continuous variables \circ If α is small, the probability that we are erroneously rejecting H_0 is very small \circ Set by us, typically at 5% \bullet Critical value z_α : \circ For two-sided: $z_{\alpha/2}$, $z_{1-\alpha/2}$ \circ For one-sided upper tail: $z_{1-\alpha}$ \circ For one-sided lower tail: z_α \circ Associated z-score with α \circ Corresponds to threshold c prior to z-score transformation \bullet P-value p : \circ For two-sided: $p = P(z | z_n \geq |z|)$ \circ For one-sided upper tail: $p = P(z \geq z_n)$ \circ For one-sided lower tail: $p = P(z \leq z_n)$ \bullet Probability, given H_0 that we observe a value as or more extreme as the observed value z_n \bullet Smallest significance level resp. largest confidence level, at which we can reject H_0 given the sample observed \circ If p-value is less than significance level α level resp. if observed value z_n is more extreme than critical value z_α resp. if observed mean \bar{x} is more extreme than threshold c , reject H_0 , because the probability that we are erroneously doing so is very small \bullet Confidence level: $1 - \alpha$, probability, given H_0 , that we retain H_0 \bullet Beta: $\beta = p$ (type II error) \bullet Power: $\circ 1 - \beta = p(\bar{x} \geq c | H_1) = p(z_n \geq z_\alpha | H_1)$ \bullet Probability, given H_A , that we reject H_0 \bullet Test statistic: $z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$
 $z|H_0 \sim N(0,1)$, $z|H_1 \sim N(\frac{\mu_1 - \mu_0}{\sigma/\sqrt{n}}, 1)$ \bullet Rejection region: $\bar{x} > \mu_0 + \sigma z_\alpha / \sqrt{n}$ \bullet Example: \circ If $X \sim N(\theta, 1)$, $H_0: \theta = 0$: $p(\bar{x} \geq c | H_0) = p(\sqrt{n} \bar{x} \geq \sqrt{n}c | H_0) = p(z_n \geq \sqrt{n}c | H_0) = 1 - \Phi(\sqrt{n}c)$ where Φ is the CDF of the normal distribution $\bullet z_n | H_0 = \frac{\bar{x} - \mu_0}{1/\sqrt{n}} = \sqrt{n} \bar{x}$ $\circ c = \frac{1}{\sqrt{n}} \Phi^{-1}(1 - \alpha)$ $\circ 1 - \beta = p(\bar{x} \geq c | H_1) = p(z_n \geq \sqrt{n}c - 1 | H_1) = p(z_n \geq \sqrt{n}(c-1) | H_1) = p(z_n \geq \sqrt{n}(-c-1) | H_1) = 1 - \Phi(\sqrt{n}(-c-1)) = 1 - \Phi(\sqrt{n}(c+1))$ where Φ is the CDF of the normal distribution $\bullet z_n | H_1 = \frac{\bar{x} - 1}{1/\sqrt{n}} = \sqrt{n}(\bar{x} - 1)$
 \bullet We can switch from $|H_1|$ to H_2 because the two distributions follow the same form, just shifted and $1 - \beta = p(z_n \geq z_\alpha | H_1) = p(z_n \geq z_\alpha - \delta | H_0) = 1 - \Phi(z_\alpha - \delta)$ where $\delta = \frac{\mu_1 - \mu_0}{\sigma/\sqrt{n}}$
Multiple comparisons problem \rightarrow Accumulation of false positive rate (α) for K tests, due to independence of tests: $P(\text{false rejections of } H_0 > 0) = 1 - P(\text{false rejections of } H_0 = 0) = 1 - (1 - \alpha)^K$
Corrections for multiple comparisons problem \rightarrow Bonferroni correction: New significance level set to $\alpha = \alpha/K$
Neyman Pearson test \rightarrow \bullet Max. power while controlling type I errors \bullet Sets α such that $\alpha \geq p$ (type I error) \bullet Then minimizes p (type II error) \bullet This is achieved by a likelihood-ratio test with threshold θ , such that α equals or is as close as possible to p (type I error): \circ If $\Lambda(x) = \frac{p(x|X)}{p(x|P_A)} > \theta$, we reject H_0 \circ Then, we have $P(\Lambda(x) > \theta | H_0) = P(\frac{p(x|P_0)}{p(x|P_A)} > \theta | H_0)$
ML Decision Rule \rightarrow If $L(x) = \log(\frac{p(x|P_A)}{p(x|P_0)}) > 0$, we reject H_0 Bayesian $M(x) = \log(\frac{P(x|P_A)P(P_A)}{P(x|P_0)P(P_0)}) > 0$, we reject H_0
Hypothesis Testing \rightarrow If $\Lambda(x) = \frac{p(x|P_A)}{p(x|P_0)} > \theta$, we reject H_0 \bullet $\theta = \frac{k(P_A, P_0)P(P_0)}{k(P_0, P_A)P(P_A)}$ \bullet In this case, θ subsumes both the prior P and the costs $k(x, x)$ Error probability \rightarrow

p (not chosen \rightarrow Information Theory Entropy \rightarrow $\bullet H(x) = -\mathbb{E}[-\log p(x)] = -\sum_x p(x) \log(p(x))$ Properties: $\bullet H_2 \geq 0$ $\bullet H$ is maximized, when x is a uniform random variable \bullet For independent variables: $H(x, y) = H(x) + H(y)$ Conditional entropy $\rightarrow H(x|y) = -\sum_{x,y} p(x,y) \log(p(x|y))$ \bullet Properties: $\bullet -\sum_{x,y} p(x,y) \log(\frac{p(x,y)}{p(x)p(y)})$ Properties: $\bullet 0 \leq H(x|y) \leq H(x)$ with eq. if x is independent with y resp. if y completely determines x \bullet Mutual information $\rightarrow I(x; y) = \mathbb{E}[\log(p(x, y)/p(x)p(y))] = \sum_{x,y} p(x, y) \log(p(x, y)/p(x)p(y))$ $\bullet I(x; y) = H(x) - H(x|y)$ Properties: $\bullet 0 \leq I(x; y) \leq H(x)$ with eq. if y completely determines x resp. if x is independent with y \bullet KL divergence $\rightarrow KL(x; y) = \mathbb{E}[\log(p(x)/q(x))] = \sum_x p(x) \log(p(x)/q(x)) = \sum_x p(x) \log(p(x)) - \sum_x p(x) \log(q(x)) = \sum_x p(x) \log(p(x)) - 1$ \bullet const. if $q(x)$ unif. $\bullet KL(x; y) = -H(p) - \sum_x p(x) \log(q(x))$ Properties: $\bullet KL(x; y) \geq 0$ Cross entropy $\rightarrow CE(p|q) = KL(p|q) + H(p) = -\sum_x p(x) \log(q(x))$ $\bullet ELBO \rightarrow ELBO(z) = \mathbb{E}[\log(p(x, z)/q(z))] = \sum_z q(z) \log(p(x, z)/q(z))$
ML Paradigms \rightarrow \bullet Parametric approach $\bullet \theta$ as fixed, unknown quantity, X as random, and known quantity \bullet Makes point eq. \bullet Focuses on max. likelihood $p(X|\theta)$ to infer posterior $p(\theta|X)$ \bullet Only requires differentiation methods \bullet High variance, but low bias \bullet MLE est. \bullet Max. log-likelihood: $\theta = \text{argmax}_{\theta} p(\theta) = \text{argmax}_{\theta} p(y_1, \dots, y_n | x_i, \theta) = \text{argmax}_{\theta} (\prod_{i=1}^n p(y_i | x_i, \theta))$ \bullet In discrete case: $\theta = \text{argmax}_{\theta} (L) = \text{argmax}_{\theta} (\sum_{i=1}^n p(y_i | x_i, \theta))$ \bullet $\text{argmax}_{\theta} \prod_{i=1}^n p_j^{N_{ij}} = \text{argmax}_{\theta} \sum_{i=1}^n N_{ij} \log(p_j)$ where $j = 1, \dots, k$ is the number of classes $\bullet N_{ij}$ county how often the outcome class j appears in y $\bullet p_j = p(y_i = j | x_i, \theta)$ \circ We can further expand to $\theta = \text{argmax}_{\theta} (L) = \text{argmax}_{\theta} \sum_{i=1}^n N_{ij} \log(p_j) = \text{argmax}_{\theta} \sum_{i=1}^n \frac{N_{ij}}{n} \log(p_j) = \text{argmax}_{\theta} \sum_{i=1}^n \frac{N_{ij}}{n} (\log(\frac{p_j}{N_{ij}/n}) + \log(N_{ij}/n)) = \text{argmax}_{\theta} \sum_{i=1}^n \frac{N_{ij}}{n} \log(\frac{p_j}{N_{ij}/n})$ \bullet This can be solved using constrained opt. with strong duality subject to $\sum_j p_j = 1$ \circ We then get $\theta_{MAP} = (N_{ij} + v) / (n + kv)$ which minimizes the KL divergence when $p_j = p_j$
Model Opt. Gradient Descent $\beta(t+1) \leftarrow \beta(t) - \eta \nabla_{\beta} L(\theta) \beta = \beta(t)$
Model Evaluation Estimator Evaluation Criteria \rightarrow \bullet Consistency: $\theta \rightarrow \theta$ as $n \rightarrow \infty$ \bullet Bias: $\mathbb{E}(\theta) - \theta$ \bullet Unbiased: $\mathbb{E}(\theta) = \theta$ \bullet Asymptotically unbiased: $\mathbb{E}[(\theta - \theta)^2] = 0$ as $n \rightarrow \infty$
Bias Variance Tradeoff Mean squared error $\mathbb{E}[(f(x) - y)^2]$ \bullet The prior is $p(\mu | \mu_0, \sigma_0^2) = \frac{1}{\sqrt{2\pi}\sigma_0^2} \exp(-\frac{(x_i - \mu)^2}{2\sigma_0^2})$ \bullet The prior is given by: $p(\mu | \mu_0, \sigma_0^2) \propto \exp(-\frac{1}{2\sigma_0^2} \sum_{i=1}^n (x_i - \mu)^2 - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2)$ \bullet Expanding the likelihood term: $\sum_{i=1}^n (x_i - \mu)^2 = \sum_{i=1}^n (x_i^2 - 2\mu x_i + \mu^2) = \sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2$ \bullet Expanding the prior term: $(\mu - \mu_0)^2 = (\mu^2 - 2\mu\mu_0 + \mu_0^2)$ \bullet This yields: $p(\mu | x, \mu_0, \sigma_0^2) \propto \exp(-\frac{1}{2} (\frac{n}{\sigma_0^2} + \frac{1}{\sigma_0^2}) \mu^2 + (\frac{\sum_{i=1}^n x_i}{\sigma_0^2} + \frac{\mu_0}{\sigma_0^2}) \mu + (\text{constant terms}))$ where the constant terms include terms that do not depend on μ , such as $\sum_{i=1}^n x_i^2$ and μ_0^2 \bullet Based on the parametric form of the Gaussian distribution, this yields $(\frac{n}{\sigma_0^2} + \frac{1}{\sigma_0^2}) \mu^2 = \frac{1}{2\sigma_n^2} \mu^2$ and $(\frac{\sum_{i=1}^n x_i}{\sigma_0^2} + \frac{\mu_0}{\sigma_0^2}) \mu = \frac{\mu_n}{\sigma_n^2} \mu$ \bullet Solving for σ_n and μ_n : $\sigma_n^2 = \frac{1}{\frac{n}{\sigma_0^2} + \frac{1}{\sigma_0^2}} = \frac{1}{\frac{n\sigma_0^2 + \sigma_0^2}{\sigma_0^4}} = \frac{\sigma_0^4}{n\sigma_0^2 + \sigma_0^2} = \frac{n\bar{x}\sigma_0^2 + \mu_0\sigma_0^2}{n\sigma_0^2 + \sigma_0^2}$ \bullet Thus, the posterior is $p(\mu | x, \mu_0, \sigma_0^2) \sim N(\mu_n, \sigma_n^2)$ \bullet Conjugate prior: Gaussian

to zero: $\theta = \text{argmax}_{\theta} p(\theta) = \text{argmax}_{\theta} \mathbb{E}[\theta|y]$ \bullet Returns single point est. \bullet Median est. \bullet Min. mean absolute error as cost func. \bullet $\mathbb{E}(\theta|y) = \theta$ \bullet The resulting est. is the median of the posterior. Proof: \circ Bayesian cost func. given by: $\mathbb{E}[|\theta - \theta|] = \int |\theta - \theta| p(\theta|y) d\theta$ \bullet The integral splits into two parts: $= \int_{-\infty}^{\theta} (\theta - \theta) p(\theta|y) d\theta + \int_{\theta}^{\infty} (\theta - \theta) p(\theta|y) d\theta$ \bullet Taking the derivative with respect to θ : $\frac{\partial \mathbb{E}[|\theta - \theta|]}{\partial \theta}$ for each term separately \bullet $\frac{\partial}{\partial \theta} \int_{-\infty}^{\theta} (\theta - \theta) p(\theta|y) d\theta = \int_{-\infty}^{\theta} p(\theta|y) d\theta - \theta p(\theta|y)$ \bullet $\frac{\partial}{\partial \theta} \int_{\theta}^{\infty} (\theta - \theta) p(\theta|y) d\theta = - \int_{\theta}^{\infty} p(\theta|y) d\theta - \theta p(\theta|y)$ \bullet Combining the two derivatives, we get: $\int_{-\infty}^{\theta} p(\theta|y) d\theta - \theta p(\theta|y) - \theta p(\theta|y) - \int_{\theta}^{\infty} p(\theta|y) d\theta$ \bullet Setting the derivative to zero: $\int_{-\infty}^{\theta} p(\theta|y) d\theta = \int_{\theta}^{\infty} p(\theta|y) d\theta$ \bullet Since the total probability is 1, this implies: $\int_{-\infty}^{\theta} p(\theta|y) d\theta = 0.5$ \bullet Returns single point est. MAP est. Max. posterior. $\theta = \text{argmax}_{\theta} p(\theta|X) \propto \text{argmax}_{\theta} p(\theta|X) p(X)$

• For any 2 points $\mathbf{x}_1, \mathbf{x}_2$ on the decision boundary $z=0$, i.e.
 $\beta \cdot \mathbf{x}_1 = 0, \quad \beta \cdot \mathbf{x}_2 = 0$ Since $\mathbf{x}_1, \mathbf{x}_2$ are on the decision boundary,
 vector \mathbf{z} can be considered a linear combination of $\mathbf{x}_1 - \mathbf{x}_2$
 ■ Combining these equations, we get $\beta \cdot (\mathbf{x}_1 - \mathbf{x}_2) = \beta \cdot \mathbf{u}z = 0$

Opt. Objective funct. — • Likelihood:

$$L(\beta) = \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}; \beta) =$$

$$\prod_{i=1}^n \sigma(z^{(i)})^{y^{(i)}} (1 - \sigma(z^{(i)}))^{1-y^{(i)}} \quad \bullet \text{ Max.}$$
 log-likelihood: $\log L(\beta) =$

$$\sum_{i=1}^n [y^{(i)} \log \sigma(z^{(i)}) + (1-y^{(i)}) \log (1 - \sigma(z^{(i)}))] =$$

$$\sum_{i=1}^n [y^{(i)} \log \frac{1}{1+e^{-z^{(i)}}} + (1-y^{(i)}) \log \frac{e^{-z^{(i)}}}{1+e^{-z^{(i)}}}] =$$

$$\sum_{i=1}^n [y^{(i)} z^{(i)} - \log(1+e^{z^{(i)}})] \quad \bullet \text{ Min. log-loss: } -\log L(\beta)$$

Opt. —
 • $\frac{\partial -\log L(\beta)}{\partial \beta} = -\sum_{i=1}^n \frac{\partial}{\partial \beta} [y^{(i)} \log \sigma(z^{(i)}) + (1-y^{(i)}) \log (1 - \sigma(z^{(i)}))] = \sum_{i=1}^n [\sigma(z^{(i)}) - y^{(i)}] \frac{\partial z^{(i)}}{\partial \beta} =$
 With $\sum_{i=1}^n [\sigma(z^{(i)}) - y^{(i)}] \mathbf{x}^{(i)}$ • Proof: • Derivative of the
 sigmoid: $\frac{\partial \sigma(z^{(i)})}{\partial z^{(i)}} = \sigma(z^{(i)}) (1 - \sigma(z^{(i)}))$ • Derivative of
 the first term: $\frac{\partial}{\partial \beta} [y^{(i)} \log \sigma(z^{(i)})] =$

$$y^{(i)} \frac{1}{\sigma(z^{(i)})} \frac{\partial \sigma(z^{(i)})}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial \beta} =$$

$$y^{(i)} \frac{1}{\sigma(z^{(i)})} \frac{1}{\sigma(z^{(i)})} \sigma(z^{(i)}) (1 - \sigma(z^{(i)})) \mathbf{x}^{(i)} =$$

$$y^{(i)} (1 - \sigma(z^{(i)})) \mathbf{x}^{(i)} = y^{(i)} \mathbf{x}^{(i)} - y^{(i)} \sigma(z^{(i)}) \mathbf{x}^{(i)}$$
 • Derivative of the second term:

$$\frac{\partial}{\partial \beta} [(1-y^{(i)}) \log (1 - \sigma(z^{(i)}))] =$$

$$(1-y^{(i)}) \frac{1}{1 - \sigma(z^{(i)})} (-1) \frac{\partial \sigma(z^{(i)})}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial \beta} =$$

$$(1-y^{(i)}) \frac{1}{1 - \sigma(z^{(i)})} (-1) \sigma(z^{(i)}) (1 - \sigma(z^{(i)})) \mathbf{x}^{(i)} =$$

$$-(1-y^{(i)}) \sigma(z^{(i)}) \mathbf{x}^{(i)} =$$

$$y^{(i)} \sigma(z^{(i)}) \mathbf{x}^{(i)} - \sigma(z^{(i)}) \mathbf{x}^{(i)}$$

Convexity of Log Loss — • Sum of convex functions is convex • Thus,
 we need to prove convexity of $\ln(1+e^{\beta \cdot \mathbf{x}^{(i)}})$ and $-y^{(i)} \beta \cdot \mathbf{x}^{(i)}$
 and then For second term: • $\mathcal{H}(\beta) = \nabla_{\beta}^2 (-y^{(i)} \beta \cdot \mathbf{x}^{(i)}) = 0$ • $\mathcal{H} \neq 0$
 given • For first term: • $\mathcal{H}(\beta) = \nabla_{\beta}^2 \ln(1+e^{\beta \cdot \mathbf{x}^{(i)}})$

$$= \nabla (\mathbf{x}^{(i)} \times \mathbf{u}'(\mathbf{x}^{(i)} - \mathbf{v}'(\mathbf{x}^{(i)})) \frac{e^{\beta \cdot \mathbf{x}^{(i)}}}{1+e^{\beta \cdot \mathbf{x}^{(i)}}})$$

$$= \frac{1}{1+e^{\beta \cdot \mathbf{x}^{(i)}}} \times \beta \cdot \mathbf{x}^{(i)} \times \mathbf{x} \mathbf{x}^T - \frac{e^{\beta \cdot \mathbf{x}^{(i)}} \times \mathbf{x}}{(1+e^{\beta \cdot \mathbf{x}^{(i)}})^2} \times \mathbf{x}^T \times e^{\beta \cdot \mathbf{x}^{(i)}}$$

$$= \frac{e^{\beta \cdot \mathbf{x}^{(i)}} \times \mathbf{x} \mathbf{x}^T (1+e^{\beta \cdot \mathbf{x}^{(i)}}) - e^{\beta \cdot \mathbf{x}^{(i)}} \times \mathbf{x} \mathbf{x}^T \times e^{\beta \cdot \mathbf{x}^{(i)}}}{(1+e^{\beta \cdot \mathbf{x}^{(i)}})^2}$$

$$= \frac{e^{\beta \cdot \mathbf{x}^{(i)}} \times \mathbf{x} \mathbf{x}^T}{(1+e^{\beta \cdot \mathbf{x}^{(i)}})^2} \quad \mathcal{H} \succ 0. \text{ Proof:}$$

$$\mathbf{a}^T \mathcal{H} \mathbf{a} = \frac{e^{\beta \cdot \mathbf{x}^{(i)}}}{(1+e^{\beta \cdot \mathbf{x}^{(i)}})^2} \mathbf{a}^T \mathbf{x} \mathbf{x}^T \mathbf{a} = \frac{e^{\beta \cdot \mathbf{x}^{(i)}}}{(1+e^{\beta \cdot \mathbf{x}^{(i)}})^2} \|\mathbf{a}^T \mathbf{x}\|^2 \geq 0$$

Regularization — • Perfectly separable data requires regularization
 • Let weights for each class k be scaled by c as $c\tilde{\beta}_k$ • Gradient of
 log-loss with respect to c is always negative, causing gradient descent
 to grow c without bound. Proof:

$$\bullet \text{ Log loss} = \sum_{i=1}^n \ln(1+e^{c\tilde{\beta}_k \cdot \mathbf{x}^{(i)}}) - y^{(i)} (c\tilde{\beta}_k \cdot \mathbf{x}^{(i)})$$

$$\bullet \nabla_c \text{ log loss} =$$

$$\sum_{i=1}^n \frac{1}{1+e^{c\tilde{\beta}_k \cdot \mathbf{x}^{(i)}}} \times e^{c\tilde{\beta}_k \cdot \mathbf{x}^{(i)}} \times \tilde{\beta}_k \cdot \mathbf{x}^{(i)} - y^{(i)} \tilde{\beta}_k \cdot \mathbf{x}^{(i)}$$

$$= \sum_{i=1}^n \tilde{\beta}_k \cdot \mathbf{x}^{(i)} \left(\frac{e^{c\tilde{\beta}_k \cdot \mathbf{x}^{(i)}}}{1+e^{c\tilde{\beta}_k \cdot \mathbf{x}^{(i)}}} - y^{(i)} \right)$$
 • Given perfect
 separation: ■ If $y^{(i)} = 1, \tilde{\beta}_k \cdot \mathbf{x}^{(i)} > 0, \frac{e^{c\tilde{\beta}_k \cdot \mathbf{x}^{(i)}}}{1+e^{c\tilde{\beta}_k \cdot \mathbf{x}^{(i)}}} - y^{(i)} < 0$
 , i.e. ■ If $y^{(i)} = 0, \tilde{\beta}_k \cdot \mathbf{x}^{(i)} < 0, \frac{e^{c\tilde{\beta}_k \cdot \mathbf{x}^{(i)}}}{1+e^{c\tilde{\beta}_k \cdot \mathbf{x}^{(i)}}} - y^{(i)} > 0$ • Thus, for
 all $i, \nabla_c \text{ log loss} < 0$ • Thus gradient descent will cause c to grow
 without bound

Multinomial Logistic Regr.
Form. • Probability of each class is estimated via the softmax funct.
 generalizes the sigmoid funct. to multiple classes):

$$P(y=k|\mathbf{x}) = \frac{e^{\mathbf{f}_j(\mathbf{x})}}{\sum_{j=1}^K e^{\mathbf{f}_j(\mathbf{x})}} = \frac{e^{\mathbf{c}_k \cdot \mathbf{x}}}{\sum_{j=1}^K e^{\mathbf{c}_j \cdot \mathbf{x}}}$$
 where T
 is the temperature and allows to smoothly interpolate between the
 differentiable softmax ($T=1$) and the non-differentiable argmax
 ($T=0$) • The predicted class is the one with the highest probability:
 $\hat{y} = \arg \max_k P(y=k|\mathbf{x})$ • Geometrically, the softmax defines
 $K-1$ linear separating hyperplanes

