

<p>EXPERIMENTAL METHODS</p> <p>treatment/indie variable what we are able to control</p> <p>outcome/dependent variable what we observe</p> <p>stratification i.e. stratify by demographic vars before splitting into ctrl groups</p> <p>double-blind even experimenters don't know who is placebo</p> <p>observational study experiment not randomizable</p> <p>control variable i.e. income level on education</p> <p>confounding variable influences both treatment and outcome variables</p>	
<p>HYPOTHESIS TESTING WITH HYPERGEOMETRIC</p> <p>aka Fisher's exact test</p> <p>does drug help? we assume H_0 is $\mu_1 = \mu_2$ so we use Hypergeometric distribution to check what is probability of drawing these numbers of deaths between study and control group under H_0.</p> <p>p-value: $\sum_{i=0}^{39} P_{H_0}(T = i)$</p> <p>pro doesn't assume knowledge of true value in ctrl group</p> <p>con assumes knowledge of margins</p>	
<p>z-test if we know the true value of σ</p> <p>t-test if we estimate $\hat{\sigma}: \frac{\bar{X}_n}{\frac{\hat{\sigma}}{\sqrt{n}}} \sim t_{n-1}$</p>	

<p>LOAD CSV DATA</p> <pre>import pandas as pd data = pd.read_csv('gamma-ray.csv')</pre>	<p>PANDAS MISC</p> <p>get first row</p> <pre>df.iloc[0]</pre>
<p>LIKELIHOOD RATIO TEST (POISSON)</p> <pre>from scipy import stats n = 1 d = 1 for row in data.data.iterrows(): G_i = row["count"] t_i = row["seconds"] lambda_i = G_i/t_i n *= stats.poisson.pmf(G_i, lambda_mle) d *= stats.poisson.pmf(G_i, lambda_i) test_statistic = -2*np.log(n/d) test_statistic</pre>	<p>QQ PLOTS</p> <p>here we were trying to figure out if we should apply a log transformation to some fields</p> <pre>stats.probplot(data[key], dist="norm", plot=plt) plt.show() stats.probplot(np.log(data[key]), dist="norm", plot=plt) plt.show()</pre>

<p>PYTHON MISC</p> <p>convert row vector into column vector</p> <pre>Y = Y.reshape(-1,1)</pre> <p>sigmoid</p> <pre>from scipy.stats import logistic logistic.cdf(x)</pre>	
<p>TODO</p> <p>Linear regression</p> <p>* multivariable</p> <p>* remove insignificant variables</p>	

LINEAR CLASSIFIERS	STOCHASTIC GRADIENT DESCENT	KERNEL PERCEPTRON	MISC
boundary line: $\theta \cdot x + \theta_0 = 0$	update: sample i at random	$\theta = 0$ for $i = 1 \dots n$ <div> if $y^{(i)} \cdot \theta \cdot \phi\left(x^{(i)}\right) \leq 0$ <div>$\theta \Leftarrow \theta + y^{(i)} \phi\left(x^{(i)}\right)$</div> </div>	$\vec{a} \cdot \vec{b} = \sum a_i b_i$ $a \cdot a = \ a\ ^2$ $\theta_z \cdot Ax = \theta_z^T Ax = \left(A^T \theta_z\right) \cdot x$
PERCEPTRON	$\theta \Leftarrow \theta - \eta_T \nabla_{\theta} \left[Loss_h\left(y^{(i)} \theta \cdot x^{(i)}\right) + \frac{\lambda}{2} \ \theta\ ^2 \right]$	$\theta \cdot \phi\left(x^{(i)}\right) = \sum_j^n \alpha_j y^{(j)} \phi\left(x^{(j)}\right) \phi\left(x^{(i)}\right)$ where α_j is nb of mistakes made on j th example	distance from point to line $d = \frac{ \theta \cdot x + \theta_0 }{\ \theta\ }$
$y^{(i)} \in \{-1, 1\}$	$= \theta \left[\begin{cases} 0 & \text{if } Loss_h = 0 \\ -y^{(i)} \cdot x^{(i)} & \text{if } Loss_h > 0 \end{cases} + \lambda \theta \right]$	we can then rewrite the perceptron algorithm without having to transform $x \Rightarrow \phi(x)$ anymore	derivative of norm squared $\left(\ h(x)\ ^2\right)' = 2h(x)h'(x)$
algorithm $\theta = \vec{0}$ for t = 1...T for i = 1...n if $y^{(i)}\left(\theta \cdot x^{(i)} + \theta_0\right) \leq 0$ then $\theta = \theta + y^{(i)} x^{(i)}$ $\theta_0 = \theta_0 + y^{(i)}$	LINEAR REGRESSION	$\theta = 0$ for $i = 1 \dots n$ <div> if $y^{(i)} \sum_j^n \alpha_j y^{(j)} K\left(x^{(j)}, x^{(i)}\right) \leq 0$ <div>$\alpha_i = \alpha_i + 1$</div> </div>	derivative $[y(\theta \cdot x)]' = yx$ because we derive for each θ_i
empirical risk $R_n(\theta) = \frac{1}{n} \sum \frac{\left(y^{(t)} - \theta x^{(t)}\right)^2}{2}$ small deviations less penalised than large ones $\nabla_{\theta} R_n(\theta) = -\left(y^{(i)} - \theta x^{(t)}\right) \cdot x^{(t)}$	gradient descent $\theta = 0$ for: random $t = \{1 \dots n\}$ # stochastic, could also loop $\theta = \theta + \eta_k \left(y^{(t)} - \theta \cdot x^{(t)}\right) \cdot x^{(t)}$	kernel composition $\cdot f: \mathbb{R}^d \rightarrow R$ and $K(x, x')$ $\widetilde{K}(x, x') = f(x)K(x, x')f(x')$ $\cdot K(x, x') = K_1(x, x') + K_2(x, x')$ $\cdot K(x, x') = K_1(x, x') + K_2(x, x')$	L1 norm $\ w\ _1 = \sum_i w_i $
REGULARIZATION	with $\eta_k = \frac{1}{1+k}$	radial basis kernel $K(x, x') = \exp\left(-\frac{1}{2}\ x - x'\ ^2\right)$	L2 norm $\ w\ _2 = \sqrt{\sum_i^n w_i^2}$
we want large boundaries: $\vdots \mid \vdots$ where \vdots are $\theta \cdot x^{(i)} + \theta_0 = \{1, -1\}$ and \mid is $\theta \cdot x^{(i)} + \theta_0 = 0$ \cdot prevents overfitting \cdot the bigger $\ \theta\ $ the faster (= shorter distance) these reach -1 and 1 \cdot we can control the margins distance with $\ \theta\ $	NONLINEAR CLASSIFICATION not linearly separable \Rightarrow map x to higher dimension	\Rightarrow correctly classifies every point \Rightarrow radial in $\phi(x)$ space	cosine similarity $\cos\left(x^{(i)}, x^{(j)}\right) = \frac{x^{(i)} \cdot x^{(j)}}{\ x^{(i)}\ \ x^{(j)}\ }$
signed distance $\gamma_i = y^{(i)} \frac{\theta \cdot x^{(i)} + \theta_0}{\ \theta\ } = \frac{1}{\ \theta\ }$ or $-\frac{1}{\ \theta\ }$ to maximise $\frac{1}{\ \theta\ }$ we minimise $\frac{1}{2} \ \theta\ ^2$	kernels i.e. $K(x, x') = \phi(x) \cdot \phi(X') = (x \cdot x') + (x \cdot x')^2$ where $\phi(x) = (x_1, x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)^T$ \Rightarrow dot product is cheap to compute	graph classification look at number of bends in figure to determine dimension	euclidean distance $dist\left(x^{(i)}, x^{(j)}\right) = \ x^{(i)} - x^{(j)}\ $
hinge loss $Loss_h(z) = Loss_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) = \max(0, 1 - y^{(i)}(\theta \cdot x^{(i)} + \theta_0))$ $= \begin{cases} 0 & \text{if } z \geq 1 \quad \text{correct side, outside margins} \\ 1 - z & \text{if } z < 1 \end{cases}$	COLLABORATIVE FILTERING assumes X has low rank $J(u, v) = \sum_{a, i \in D} \frac{(Y_{ai} - u_a v_i)^2}{2} + \frac{\lambda}{2} \sum u_a^2 + \frac{\lambda}{2} \sum v_i^2$ $X = \begin{bmatrix} \text{user1rating1} & \text{user1rating2} \\ \text{user2rating1} & \text{user2rating2} \end{bmatrix} = u \cdot v^T$	K-NEAREST NEIGHBOURS $\hat{Y}_{ai} = \frac{1}{K} \sum_{b \in \text{KNN}(a, i)} Y_{bi}$ or weigh by similarity $\hat{Y}_{ai} = \frac{\sum_{b \in \text{KNN}(a, i)} sim(a, b) Y_{bi}}{\sum_{b \in \text{KNN}} sim(a, b)}$	
objective $J(\theta, \theta_0) = \frac{1}{n} \sum Loss_h\left(y^{(i)}\left(\theta \cdot x^{(i)} + \theta_0\right)\right) + \frac{\lambda}{2} \ \theta\ ^2$ where λ is regularisation loss and minimising that term increases margins \rightarrow by minimising J we balance maximising margins and minimising loss	\Rightarrow u: user's rating tendency \Rightarrow v: info about movie	euclidian distance $\ x_a - x_b\ $ cosine similarity $\cos(\theta) = \frac{x_a \cdot x_b}{\ x_a\ \ x_b\ }$	
gradient descent $\theta \Leftarrow \theta - \eta \nabla J(\theta)$	alternating projections set v, optimise u, set u... continue until local convergence (exercise link)	limitations my combination of tastes maybe uncommon	
SOLVING QUADRATIC		ACTIVATION FUNCTIONS	
if linearly separable and don't allow any errors, can solve quadratic: $\Rightarrow J(\theta, \theta_0) = \frac{1}{2} \ \theta\ ^2$		ReLU $f(z) = \max(0, z)$ tanh $f(z) = \tanh(z)$ linear $f(z) = z$	

DEEP LEARNING

- intermediate layer is like feature representation: "learning ϕ "
- activation layer acts like linear classifier

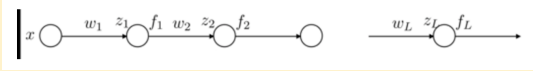
backpropagation

compute partial derivatives using chain rule to do gradient descent

$$w_1 \leftarrow w_1 - \eta \left(\frac{\delta \text{Loss}}{\delta w_1} \right)$$

$$\frac{\delta \text{Loss}}{\delta w_1} = \frac{\delta \text{Loss}}{\delta f_1} \cdot \frac{\delta f_1}{\delta z_1} \cdot \frac{\delta z_1}{\delta w_1}$$

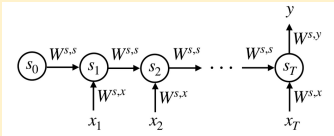
risks are vanishing or exploding gradients with long chains



RNN

- same params for each layer
- variable number of layers
- ⇒ encodes sentences

$$s_t = \tanh(W^{s,s} s_{t-1} + W^{s,x} x_t)$$



RNN have some gradient issues

gating network

in default RNN we overwrite s_t but sometimes it's useful to retain control over what's overwritten at each step

$$g_t = \text{sigmoid}(W^{g,s} s_{t-1} + W^{g,x} x_t)$$
$$s_t = (1 - g_t) \odot s_{t-1} + g_t \odot \tanh(W^{s,s} s_{t-1} + W^{s,x} x_t)$$

where \odot is element-wise multiplication

$$(1 - g_t) : \begin{pmatrix} 0 \\ 1 \end{pmatrix} \rightarrow \begin{pmatrix} 1 \\ 0 \end{pmatrix} \text{ i.e. } g_t \text{ is use to control what we retain from previous layer}$$

MARKOV MODELS

<beg> w_1 w_2 w_3 <end>
 w_0 w_L

$$P(w_1, ..., w_L) = P(w_1 | w_0) ... P(w_L | w_{L-1})$$

$$\text{with } P(w' | w) = \frac{\text{count}(w, w')}{\sum_{\tilde{w}} \text{count}(w, \tilde{w})}$$

RNNs for sequences

variable history

$$\phi(\text{word}) : X^W : s_t X : p$$
$$s_{t-1} : X$$

i.e. we're feeding the previous state as input to the next

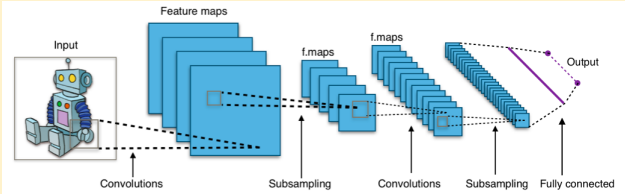
CNN

convolution

- classify by small patches
- slide over input
- use the same weights for each patch to build feature map

pooling (subsampling)

- stores if there's an activation for feature in that area
- stride



cross-correlation/convolution

we don't actually need to flip $g(t)$ as we would with convolution, because we're learning the filter, the values are just flipped

$$1d \ f \odot g(t) = \sum f(T)g(T+t)$$

$$md \ f \odot g(x, y) = \sum_{T_1} \sum_{T_2} f(T_1, T_2)g(T_1 + x, T_2 + y)$$

K-MEANS

- 1) randomly assign $z^{(1)} ... z^{(2)}$
- 2) iterate
 - a) assign x to closest z
$$\text{Cost} \left(z^{(1)} ... z^{(k)} \right) = \sum_{i=1}^n \min_{j=1...k} \left\| x^{(i)} - z^{(j)} \right\|^2$$
 - b) given $C_1 ... C_k$ find best rep. z in each C
$$\text{Cost}(C_1 ... C_k) = \min_{z^{(1)} ... z^{(k)}} \sum_{j=1}^k \sum_{i \in C_j} \left\| x^{(i)} - z^{(j)} \right\|^2$$

problems with K-means

- points are not in set of X s
- (euclidean distance)² only
- K-means is $O(n \cdot k \cdot d)$ for one update step (2)

find best z_j for C_j

$$\frac{\delta}{\delta z_j} \sum_{i \in C_j} \left\| x^{(i)} - z^{(j)} \right\|^2 = 0$$

$$z^{(j)} = \frac{\sum x^{(i)}}{|C_j|}$$

where $|C_j|$ is size of C_j

impact of init

- init matters
- we can end up with undesirable convergence

K-MEDIANS

Like K-means but with medians. Compute median in each single dimension in the Manhattan distance formulation (see [ex.](#))

K-MEDOIDS

- 1) pick init $z^{(i)} ... z^{(j)}$ from points X
- 2) iterate until no change in cost
 - a) for $i = 1 ... n$

$$C_j = \left\{ i : z^{(j)} \text{ is closest to } x^{(i)} \right\}$$

- b) pick the most central $x^{(i)}$ in C_j to be $z^{(j)}$

for $j = 1 ... k$

$$z_j = \left\{ x^{(i)} ... x^{(n)} : \sum_{i \in C_j} \text{dist} \left(x^{(i)}, z^{(j)} \right) \text{ is minimal} \right\}$$

· same idea as K-means but z_j is in x and can use any distance measure

· K-medoids is $O(n^2 \cdot k \cdot d)$ for update step (2)

Δ step 2 gets executed even if no change and cost, and order we cycle in matters and can update the medoid (see [ex.](#))

GENERATIVE MODELS

$p(w | \theta) = \theta w$: probability of generating word w given all possibilities`

probability of generating doc D

$$p(D | \theta) = \prod_{i=1}^n \theta_{w_i} = \prod_{w \in W} \theta_w^{\text{count}(w)}$$

MLE

$$\max P(D | \theta) \rightarrow \sum_{w \in W} \text{count}(w) \log(\theta_w)$$

$$\Rightarrow \tilde{\theta}_w = \frac{\text{count}(w)}{\sum_{w' \in W} \text{count}(w')}$$

prediction

essentially acts as a linear classifier through origin

$$\log \frac{p(D | \theta^+)}{p(D | \theta^-)} = \begin{cases} \geq 0 & + \\ < 0 & - \end{cases}$$

prior, posterior, likelihood

for when we have same prior knowledge of document classification

$$P(y = + | D) = \frac{P(D | \theta^+) P(y = +)}{P(D)}$$

⇒ linear classifier with offset influenced by prior

Δ N-1 params in multinomial (i.e. remember the binomial uses 1 param for two values)

GAUSSIAN GENERATIVE MODELS	MARKOV DECISION PROCESS	RL
<p>$X \mid \mu, \sigma^2 \sim N(\mu, \sigma^2)$</p> <p>want to find μ and σ^2 that give highest likelihood to training data.</p> <p>$\Rightarrow \frac{\delta l_n}{\delta \mu} = 0 \text{ and } \frac{\delta l_n}{\delta \sigma^2} = 0$</p> <p>$\widehat{\mu} = \frac{1}{n} \sum x^{(i)} \text{ and } \widehat{\sigma^2} = \frac{1}{nd} \sum \ x^{(i)} - \mu\ ^2$</p>	<p>· states $s \in S$</p> <p>· actions $a \in A$</p> <p>· action dependent on transition probabilities</p> <p>$T(s, a, s') = P(s' \mid a, s)$</p> <p>$\sum_{s' \in S} T(s, a, s') = 1$</p> <p>$R(s, a, s')$: reward for starting s, doing a, ending s'</p>	<p>typically won't have T and R ahead of time</p> <p>Q-value iteration</p> <p>want to calculate estimates for Q but don't have T and R. we'll use exponential running average to get estimates for $Q(s, a)$:</p> <p>$Q_{i+1}(s, a) = \alpha \cdot \text{sample} + (1 - \alpha)Q_i(s, a)$</p> <p>sample: $R(s, a, s') + \gamma \max_{a'} Q_i(s', a')$</p>
<p>MIXTURE MODELS</p> <p>mixture components K: $N(x, \mu_j, \sigma_j^2)$</p> <p>mixture weights $p_1 \dots p_k \sum p_j = 1$</p> <p>$\theta = p_1 \dots p_k, \mu_1 \dots \mu_k, \sigma_1^2 \dots \sigma_k^2$</p> <p>$p(x \mid \theta) = \sum_{j=1}^k p_j N(x, \mu_j, \sigma_j^2)$</p> <p>$p(S_n \mid \theta) = \prod_{i=1}^n \prod_{j=1}^k p_j N(x_i, \mu_j, \sigma_j^2)$</p>	<p>UTILITY FUNCTIONS</p> <p>finite horizon (not important)</p> <p>count reward for next n steps: $U(s_0 \dots s_n) = \sum^n R(s_i)$</p> <p>infinite horizon/discounted reward</p> <p>$U(s_0 \dots s_n) = \sum_{k=0}^\infty \gamma^k R(s_k) \leq \frac{R_{\max}}{1 - \gamma}$</p> <p>value function</p> <p>$V^*(s)$: expected reward if agent acts optimally starting at s</p> <p>$\pi: S \rightarrow A$: assigns action $\pi(s)$ to state s</p> <p>optimal policy π^* assigns action at every state that maximises expected utility</p>	<p>1) initialisation $Q(s, a) = 0 \forall a$</p> <p>2) iterate until convergence</p> <p>a) collect samples: $s, a, s', R(s, a, s')$</p> <p>b) $Q_{i+1}(s, a) = \alpha \left[R(s, a, s') + \gamma \max_a Q_i(s', a') \right] + (1 - \alpha)Q_i(s, a)$</p> <p>$\epsilon$ greedy</p> <p>balance exploration and exploitation</p> <p>· random action w.p. ϵ</p> <p>· best current action w.p. $(1 - \epsilon)$</p> <p>ϵ should decay over time</p>
<p>EM algorithm</p> <p>"random" init of $\theta: \mu^{(i)}, \sigma^{(i)}, p_i$</p> <p>init matters for speed and outcome</p> <p>· converges locally</p> <p>· good idea to use k-means for init</p> <p>expectation step</p> <p>$p(j \mid i) = \frac{p_j N(x_i, \mu_j, \sigma_j^2 I)}{p(x \mid \theta)}$</p> <p>$\Rightarrow$ probability that point i in cluster j</p> <p>$\Rightarrow p(x \mid \theta) = \sum_{j=1}^k p_j N(x_i, \mu_j, \sigma_j^2 I)$</p> <p>maximisation step</p> <p>re-estimate parameters, (nb: cluster j, point i)</p> <p>$\widehat{n}_j = \sum_{i=1}^n p(j \mid i)$</p> <p>$\widehat{p}_j = \frac{\widehat{n}_j}{n}$</p> <p>$\widehat{\mu}_j = \frac{1}{\widehat{n}_j} \sum_{i=1}^n p(j \mid i) x^{(i)}$</p> <p>$\widehat{\sigma}_j^2 = \frac{1}{\widehat{n}_j d} \sum_{i=1}^n p(j \mid i) \ x^{(i)} - \mu^{(j)}\ ^2$</p>	<p>$Q^*(s, a)$: expected reward starting at s, taking action a and acting optimally</p> <p>BELLMAN EQUATIONS</p> <p>$V^*(s) = \max_a Q^*(s, a) = Q^*(s, \pi^*(s))$</p> <p>$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$</p> <p>VALUE ITERATION ALGORITHM</p> <p>$V_k^*(s) \rightarrow_{k \rightarrow \infty} V(s)$</p> <p>· init: $V_0^*(s) = 0$</p> <p>· iterate until $V_k^*(s) \approx V_{k+1}^*(s) \forall s$</p> <p>$V_{k+1}^* = \max_a \left[\sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k^*(s')] \right]$</p> <p>· compute $Q^*(s, a)$</p> <p>then $\pi^*(s) = \arg \max_a Q^*(s, a)$</p> <p>Q-VALUE ITERATION</p> <p>same as above but with Q and the update is:</p> <p>$Q_{k+1}^* = \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma \max_{a'} Q_k^*(s', a') \right]$</p> <p>$\Rightarrow$ we're summing over all possible states we may land in</p>	