

CHEM 313 "Midterm"

General

$$\mathcal{X} = \mathbb{R}^D \quad Y = \text{prediction}$$

$$N = \# \text{ samples} \quad T = \text{test set} \quad \sigma^2(a) = 1 - \sigma(a)$$

$$\sigma(a) = (1 + e^a)^{-1} \quad \sigma'(a) = \sigma(a)(1 - \sigma(a))$$

$$\text{tanh}(a) = 2\sigma(a) - 1 \quad \text{tanh}'(a) = 1 - \text{tanh}^2(a)$$

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)} \quad \text{posterior} \quad \text{evidence}$$

Probability Transforms

$$Y = f(X) \quad X \sim p_X \quad \int p_X(x) dx = 1$$

$$X: \mathcal{X} \rightarrow \mathcal{Y} \quad Y: \mathcal{Y} \rightarrow \mathcal{Z}$$

$$\bar{X} = \text{Base object} \quad \bar{Y} = \text{Base object}$$

$$f: (\mathcal{X}, \bar{X}) \rightarrow (\mathcal{Y}, \bar{Y})$$

$$Y: \mathcal{Y} \rightarrow \mathcal{Z} = Y(\omega) = f(X(\omega))$$

$$p_Y(B) = P(Y \in B) = P(X \in f^{-1}(B)) = p_X(f^{-1}(B))$$

$$p_Y(y) = \int p_X(x) \delta(y - f(x)) dx$$

$$= p_X(f^{-1}(y)) \left| \frac{dy}{dx} (f^{-1}(y)) \right| \quad \text{Monotonic}$$

$$= \sum_{\{x | f(x)=y\}} \frac{p_X(x)}{|f'(x)|} \quad \text{General (IFT)}$$

Dense

$$S(f(x)) = \sum_{\{x_i | f(x_i) = y\}} \frac{S(x_i)}{|f'(x_i)|}$$

Normal

$$\mu_{A|B} = \mu_A + \sum_{AB} \Sigma_{AB}^{-1} (\mu_B - \mu_B)$$

$$\Sigma_{A|B} = \Sigma_{AA} - \sum_{AB} \Sigma_{AB}^{-1} \Sigma_{BB} \Sigma_{BA}$$

$$\log(p(x)) = -\frac{1}{2} \left(\ln(2\pi) + (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

Conditional Expectation

$$E[Y|X] = \int y p_{Y|X}(x, y) dy$$

$$p_{Y|X}(x, y) = \frac{p_{XY}(x, y)}{p_X(x)}$$

$$\frac{p(y|x, w, \beta)}{\text{likelihood}} = \frac{p(y|x, w, \beta)}{\text{prior}}$$

Multiple Outputs

$$y \in \mathbb{R}^K$$

$$x \in \mathcal{X}$$

$$y_{nk} = w_{nk} \phi(x_n)$$

$$\ln(p(y|X, w, \beta)) = \sum_n \ln(N \text{tr}(w_{nk} \phi_{nk}^T \Sigma^{-1} \phi_{nk}))$$

$$= \frac{NK}{2} \ln\left(\frac{\beta}{2\pi}\right) - \frac{\beta}{2} \sum_{n,k} \phi_{nk}^T (y_{nk} - w_{nk})^2$$

$$w_{nk} = (\phi_{nk} \phi_{nk}^T)^{-1} y_{nk} + w_{nk} \quad (K \text{ rows, all indep})$$

Regularization

$$E_D(w) + \lambda E_W(w) \quad E_W(w) = \frac{1}{2} \sum |w_i|^d$$

"Ridge" $\alpha=2$: $w = (\lambda I + \phi^T \phi)^{-1} \phi^T y$

Loss: $J = \sum_n \left(\frac{\beta}{2} \max\left(\left| \frac{\beta}{2} - \frac{1}{2}, 0 \right| \right) \right)$

Loss: $J = \sum_n \left(\frac{\beta}{2} \max\left(\left| \frac{\beta}{2} - \frac{1}{2}, 0 \right| \right) \right)$

Loss: $J = \sum_n \left(\frac{\beta}{2} \max\left(\left| \frac{\beta}{2} - \frac{1}{2}, 0 \right| \right) \right)$

Linear

$$x_i \in \mathcal{X} \quad y \in \mathbb{R}$$

$$\phi: \mathcal{X} \rightarrow \mathbb{R}^{M-1}$$

$$y(x, w) = w_0 + \sum_{i=1}^{M-1} w_i \phi_i(x)$$

$$= w^T \phi(x)$$

ML Linear

$$t = y + \epsilon$$

$$\epsilon \sim N(0, \beta^{-1}) \quad \text{assumes}$$

$$p_{t|y}(t, y) = N(t - y, \beta^{-1})$$

$$p_{t|x}(t, x) = N(t | y(x, w), \beta^{-1})$$

$$E[t|x] = \int t p(t|x) dt = y(x, w)$$

assumes iid samples $y_n = \phi(x_n) w$

$$p(\tilde{t} | X, w, \beta) = \prod_{n=1}^N N(t_n | y_n, \beta^{-1}) \quad \text{Likelihood}$$

$$\lambda = \ln(p(t | X, w, \beta)) = \frac{N}{2} \ln(\beta) - \frac{N}{2} \ln(2\pi) - \beta E_D(w)$$

$$E_D(w) = \frac{1}{2} \sum_{n=1}^N (t_n - y_n)^2$$

$$\frac{\partial \lambda}{\partial w_k} = 0 = \sum_n (t_n - w_k \phi_{nk}) \frac{\partial}{\partial w_k} (-w_k \phi_{nk})$$

$$= - \sum_n (t_n - w_k \phi_{nk}) \phi_{nk}$$

$$\phi_{nk} \phi_{nk}^T = \phi_{nk}^T \phi_{nk}$$

$$\phi^T \phi = \phi^T \phi$$

$$\tilde{w} = (\phi^T \phi)^{-1} \phi^T t$$

$$\text{Bias: } w_0 = \frac{1}{N} \sum_{n=1}^N y_n \phi_0(x_n)$$

ie offset diff

Regularization

$$E_D(w) + \lambda E_W(w) \quad E_W(w) = \frac{1}{2} \sum |w_i|^d$$

"Ridge" $\alpha=2$: $w = (\lambda I + \phi^T \phi)^{-1} \phi^T y$

Loss: $J = \sum_n \left(\frac{\beta}{2} \max\left(\left| \frac{\beta}{2} - \frac{1}{2}, 0 \right| \right) \right)$

Loss: $J = \sum_n \left(\frac{\beta}{2} \max\left(\left| \frac{\beta}{2} - \frac{1}{2}, 0 \right| \right) \right)$

Loss: $J = \sum_n \left(\frac{\beta}{2} \max\left(\left| \frac{\beta}{2} - \frac{1}{2}, 0 \right| \right) \right)$

Multiple Outputs

$$y \in \mathbb{R}^K$$

$$x \in \mathcal{X}$$

$$y_{nk} = w_{nk} \phi(x_n)$$

$$\ln(p(y|X, w, \beta)) = \sum_n \ln(N \text{tr}(w_{nk} \phi_{nk}^T \Sigma^{-1} \phi_{nk}))$$

$$= \frac{NK}{2} \ln\left(\frac{\beta}{2\pi}\right) - \frac{\beta}{2} \sum_{n,k} \phi_{nk}^T (y_{nk} - w_{nk})^2$$

$$w_{nk} = (\phi_{nk} \phi_{nk}^T)^{-1} y_{nk} + w_{nk} \quad (K \text{ rows, all indep})$$

Classification Metrics

Confusion matrix:

	True (+)	False (-)
A =	True (+)	False (-)

$$Accuracy = \frac{A_{ii}}{\sum_j A_{ij}}$$

$$E[L] = \sum_i \sum_j \int_{\mathcal{R}_j} dx p(x, C_i)$$

\mathcal{R}_j decision regions

$$p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)}$$

$$p(C_k|x, C_j) \propto p(x|C_k)p(C_k)$$

2 features \rightarrow "Naive Bayes" conditional indep assumption

used to classify \rightarrow $p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(C_k)}$ Bayes' theorem

Cross Entropy

$$H(p, q) = -\sum_{n=1}^N \frac{1}{N} (t_n \ln(t_n) - (1 - t_n) \ln(1 - t_n))$$

Regression Metrics

$$E[L] = \int \int dx dt L(t, y(x)) p(x, t) dx dt$$

usually

$$= \int \int dx dt (t - y)^2 p(x, t) dx dt$$

$$\frac{\delta E}{\delta y} = 2 \int (y(x) - t) p(x, t) dx = 0$$

$$y(x) \int p(x, t) dt = \int t p(x, t) dt$$

$$y(x) p(x) = \int t p(x, t) dt$$

$$y(x) = E[t|x] \quad \text{functionally optimized solution to loss}$$

$$(y(x) - t)^2 = (y(x) - E[t|x])^2 + 2(y(x) - E[t|x])(E[t|x] - t) + (E[t|x] - t)^2$$

0 in integral

$$E[L] = \int (y(x) - E[t|x])^2 p(x) dx + \int \text{bias}$$

$$\int (E[t|x] - t)^2 p(x) dx \quad \text{variance (irreducible minimum)}$$

"rigid" = few params \rightarrow high bias \rightarrow low variance \rightarrow (high λ)

"flexible" = many params \rightarrow low bias \rightarrow high variance \rightarrow (low λ)

Bayesian Lin Reg

$$p(w) = N(w | \mu_w, \Sigma_w) \text{ prior}$$

$$p(w|t) = N(w | \mu_w, \Sigma_w) \quad y_n = w^T \phi(x_n)$$

$$\mu_{w_i} = \sum_j S_{ij} \phi_j^T \mu_{y_n} + \beta \phi_n^T t_n$$

$$\Sigma_{w_i}^{-1} = \Sigma_{w_i}^{-1} + \beta \phi_n \phi_n^T \quad \Sigma_0^{-1} = \alpha I$$

$$w_{MAP} = \mu_{w_i}$$

Ch4 Classification

$$y(x) = w_i x_i$$

$$y: y \geq 0 \Rightarrow C_0 \\ \text{else} \Rightarrow C_1$$

$$K: \begin{pmatrix} K \\ 2 \end{pmatrix} \text{ classifiers to distinguish pairs} \quad \text{sign}(y) = \text{classifer}$$

$$w_i x_i = 0 \text{ hyperplane divides}$$

$$\text{1 of } K \text{ Least Squares} \\ + \epsilon [0, 1]^K \quad K_2 = \arg \max(t)$$

$$y_k = w_i x_i$$

$$w_{ik}^P = (\phi_{i1} \phi_{ij})^T \phi_{j1} t_{ik}$$

$$\rightarrow \text{if all } t \text{ s.t. } \sum t_k = 1, \text{ all predictions will also sum to 1}$$

$$\rightarrow \text{very bad w outliers}$$

Fisher LDA

Intuition: maximize mean separation, minimize in-class variance

ex 2 classes:

$$\mu_{k1} = \frac{1}{N_k} \sum_{n \in N_k} x_{n1}$$

$$\mu_{k2} - \mu_{k1} = w^T (\mu_{k2} - \mu_{k1})$$

$$w^T w = 1 \\ \downarrow \\ \text{orthogonalization}$$

$$S_k^2 = \sum_{n \in N_k} (w_i x_{ni} - \mu_{k1})^2$$

$$K=2 \quad J(w) = \frac{(\mu_{k2} - \mu_{k1})^2}{S_{k1}^2 + S_{k2}^2} \\ \xrightarrow{\text{max}} = \frac{w^T S_{ij}^B w}{w^T S_{ij}^B w}$$

$$S_{ij}^B = (\mu_{k1} - \mu_{k2})(\mu_{k1} - \mu_{k2}) \\ S_{ij}^B = \sum_{n \in C_1} (x_{n1} - \mu_{k1})(x_{nj} - \mu_{kj}) + \sum_{n \in C_2} (x_{n1} - \mu_{k2})(x_{nj} - \mu_{kj})$$

$$\text{Optimal: } S_{ij}^B w_j \propto \mu_{k1} - \mu_{k2}$$

$$w_i \propto (S_{ij}^B)^{-1} (\mu_{k1} - \mu_{k2})$$

$$\Leftrightarrow \text{Least squares} \\ \text{where for } C_1, t_n = \frac{1}{N_1} \\ C_2: t_n = \frac{-N}{N_2}$$

Perceptron

$$\Theta: \mathbb{R} \rightarrow \{-1, 1\} \text{ (Heaviside)}$$

$$y(x) = \Theta(w_i \phi_i(x) + b) \quad \left[\begin{array}{l} \text{maximally} \\ \text{nonlinear} \end{array} \right]$$

$$L = \sum_{n \in \text{misclassified}} 1$$

$$L_n = - \sum_i w_i \phi_i(x_n) t_n \geq 0$$

Probabilistic Classification Data \rightarrow Sigmoid predictor with Gaussian

$$\text{Assume } p(x|C_k) = \frac{1}{(2\pi)^{D/2} |K|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu_k)^T \Sigma_{ij}^{-1} (x_j - \mu_{kj})\right)$$

$$p(C_k|x) = \sigma\left(\underbrace{w_{ik} \phi_i(x) + w_{0k}}_{\theta_{ik}(x)}\right)$$

$$w_{ik} = \sum_j \phi_j^T y_{kj}$$

$$w_{0k} = -\frac{1}{2} \mu_k^T \Sigma_{ij}^{-1} \mu_{kj} + \ln(p(C_k))$$

Max-Likelihood

$$t_n = 1 \Rightarrow C_1 \quad p(C_1) = \pi \\ t_n = 0 \Rightarrow C_2 \quad p(C_2) = 1 - \pi$$

$$p(x_n, C_1) = p(C_1) p(x_n | C_1)$$

$$= \pi N(x_n | \mu_1, \Sigma^1)$$

$$p(x_n, C_2) = (1 - \pi) N(x_n | \mu_2, \Sigma^2)$$

$$p(t | \pi, \mu_1, \mu_2, \Sigma^1) = \prod_{n=1}^N \left(\pi N(x_n | \mu_1, \Sigma^1) \right)^{t_n} \left((1 - \pi) N(x_n | \mu_2, \Sigma^2) \right)^{1-t_n}$$

: ML

$$\pi = \frac{N_1}{N_1 + N_2}$$

$$\mu_{k1} = \frac{1}{N_k} \sum_{n \in N_k} x_{n1}$$

$$\Sigma^1 = \frac{1}{N} \sum_k \Sigma_k$$

$$S_k = \sum_{n \in C_k} (x_{n1} - \mu_{k1})(x_{nj} - \mu_{kj})$$

Logistic Regression

$$p(C_1 | \phi) = y(\phi(x)) = \sigma(w_i \phi_i(x))$$

$$p(t|w) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}$$

$$t_n \in \{0, 1\}$$

$$E(w) = -\ln p(t|w)$$

$$= - \sum_n \left(t_n \ln y_n + (1 - t_n) \ln (1 - y_n) \right)$$

$$\frac{\partial E}{\partial w_j} = \sum_n (y_n - t_n) \phi_j(x_n)$$

Procedure

\rightarrow Define inductive bias dist

\rightarrow Define loss

\rightarrow Weights = $-\log(L(\theta(x)))$ via some gradient descent

Tree CHM
 $Y_{\text{mean}}(x) = \frac{1}{M} \sum_{m \in \text{min}} Y_m(x)$
Bagging
 → M "bootstrap" datasets are resampled from the original to generate separate datasets / models which are aggregated via voting
 → mainly to reduce variance (and overfitting) as side effect

Boost
 $Y_m(x) = h(x) + \epsilon_m(x)$
 $E[(Y_m - h)^2] = E[\epsilon_m(x)^2]$
 $E_{\text{AV}} = \frac{1}{M} \sum_m E[\epsilon_m(x)^2]$ individually
 $E_{\text{com}} = E\left[\left(\frac{1}{M} \sum_m Y_m - h(x)\right)^2\right]$
 $= E\left[\left(\frac{1}{M} \sum_m \epsilon_m\right)^2\right]$

$E_{\text{com}} = \frac{1}{M} E_{\text{AV}}$ Assumes $E[\epsilon_m(x)] = 0$
 $E[\epsilon_m(x)\epsilon_n(x)] = \delta_{mn}\sigma^2$

Boosting
 → weak learners trained sequentially where if a previous model misclassifies, point weight gets boosted
 → weighted sum @ end
 $Y_m = \text{sign}\left(\sum_{n=1}^m Y_n(x)\right)$ $t_n \in \{-1, 1\}$

Ada Boost
 $\epsilon_m = \frac{\sum_{i=1}^n w_i I(y_i \neq (h_m(x) + t_n))}{\sum_{i=1}^n w_i}$
 $d_n = \ln\left(\frac{1 - \epsilon_m}{\epsilon_m}\right)$
 $w_n^{\text{new}} = w_n^{\text{old}} \exp(d_n I(y_i \neq (h_n(x) + t_n)))$ by boost missed ones

Ensemble: exponential
Boosting v. Bagging

trained in sequence	trained indep
imposes bias	imposes variance

Tree Models (CART)
 → binary tree partition feature space
 → each branch plays the dim & cutoff via exhaustive search greedily
 $Y_c = \frac{1}{M_c} \sum_{x \in \text{leaf } c} Y_x$ optimal prediction
 $\# \text{pts in } R_c$
 $Q_c(T) = \sum_{x \in R_c} (Y_c - Y_x)^2$
 $C(T) = \sum_{c=1}^M Q_c(T) + \lambda |T|$
 pruning: $T \Rightarrow$ too complex / bad

Losses
 $Q_c(T) = \sum_{k=1}^K p_{ck} \ln(p_{ck})$ Cross entropy
 $= \sum_{k=1}^K p_{ck} (1 - p_{ck})$ Gini index
 p_{ck} = proportion of pts in region c in class k
 $p_{ck} = 0, 1 \Rightarrow$ good model
 $p_{ck} = 1/2 \Rightarrow$ bad
 → differentiable unlike misclassification rate

SVD / PCA
 $A = U S V^T$ $U^T U = I$ $V^T V = I$ $S \geq 0$
 $U = \text{eig}(AA^T)$ cols = "left sing. vec"
 $V = \text{eig}(A^T A)$ cols = "right sing. vec"
 $\|A - A_K\|_F = \sqrt{\sum_{k=K+1}^r \sigma_k^2}$
 $r = \text{rank}(A)$ $\max_{W \in \mathbb{R}^{n \times n}} |V(W)| = W^T S W$
 $W = [w_1, \dots, w_n]$ ie top s -vals capth more variance
 $\sum w_i w_i^T = P_{CA}$ direction (ortho)

Manifold Techniques
Multidimensional Scaling
 → assumes high dim data actually on a lowly smooth low dim manifold
 → some assume global connectivity
 → isomap unimodal
 → SNE good for misn / clustering while others bad
 $C_{ij} = \delta_{ij} - \frac{1}{N} \sum_{l \neq j} C_{il} C_{lj}$ "contingency matrix" (ie SIs)
 $X'_{in} = C_{ij} X_{jn}$
 $K_{mn} = X'_{in} X'_{jm}$
 goal: find $Z_{x,n}$ s.t. pairwise distances are similar to specified D_{ij} by minimum strain

$d(z) = \sum_{n,m} (K_{mn} - Z_{n,z} Z_{m,z})^2$
 $\hat{Z}_{in} = C_{ij} Z_{jn}$
(t)SNE
 $P_{ij} = \frac{\exp(-\frac{1}{2\sigma^2} \|x_i - x_j\|^2)}{\sum_k \exp(-\frac{1}{2\sigma^2} \|x_i - x_k\|^2)}$
 prob j and "i" are neighbors in full space
 $q_{ji} = \frac{\exp(-\frac{1}{2} \|z_i - z_j\|^2)}{\sum_k \exp(-\frac{1}{2} \|z_i - z_k\|^2)}$ low dim
 z ideally
 $\alpha = \sum_{ij} P_{ij} \log \frac{P_{ij}}{q_{ji}}$
 → pulls far points together rather than per arg
 sum exp(...) for $(\|x_i - x_j\|^2)^{-1}$ int SNE

Clustering
 $N_{ij} = \# \text{ objects in class } i \text{ of class } j$
 $N_i = \sum_j N_{ij}$ $P_{ij} = \frac{N_{ij}}{N_i}$
 $p_i = \max_j P_{ij}$
 $\text{purity} = \sum_i \frac{N_i}{N} p_i$

Hierarchical Agglomerative Clustering
 $D_{ij} \geq 0$ input matrix → tree which groups its elements
 successively group most similar then goes up hierarchy until all 1 group
 $d_{ij} = \sqrt{\sum_k (x_{ik} - x_{jk})^2}$ or some other metric

K-Means
 $K, \mu_c \in \mathbb{R}^D$ minimizes within cluster variance
 $Z_n^* = \arg \min_{x \in \mathbb{R}^D} \sum_{n \in C_n} \|x_n - \mu_n\|^2$
 $\mu_{nj} = \frac{1}{N_n} \sum_{n \in C_n} x_{nj}$
 $J(M; Z) = \sum_n \left(\|x_n - Z_{n, M_n}\|^2 \right)$ distribution, opt M w/ alt minimization
 $Z \in \{0, 1\}^{N \times K}$
 $x \in \mathbb{R}^{ND}$ data
 $M \in \mathbb{R}^{DK}$ k means as columns

Vector Quantization
 $\hat{Z} = \frac{1}{N} \sum_n \|x_n - \text{decode}(\text{encode}(x_n))\|^2$
 Quantization: $\text{encode}(x_n) = \arg \min_i \|x_n - \mu_i\|^2$
 $\text{decode}(z) = \mu_z$
 Initialization: each vector gets replaced with symbol in codebook
 $x_n \in \mathbb{R}^D \rightarrow z_n \in \{1\}^K$ (snaps them)
 $O(NDB) \rightarrow O(N|N|K + KDB)$
 ↑ \uparrow
 bit scalar \uparrow inverter & codebook encoded
 K means means codebook
 VQ uses codebook

K-means vs GMM
 hard soft (probabilistic assignment)
 spherical, equally sized Gaussian (has own center)
 within cluster variance minimized maximize log likelihood prediction
 $\min_{\mu, \sigma} \sum_{n \in C_n} \|x_n - \mu_{C_n}\|^2$
 $\log p(x) = \sum_i \log \left(\sum_{n \in C_n} \frac{1}{N} \exp(-\frac{1}{2\sigma_n^2} \|x - \mu_n\|^2) \right)$

NNs
 $x = \text{layer}$ nonlinear activation
 $z_k^{\text{net}} = h\left(\sum_j w_{kj} z_j^{\text{net}}\right) = h^*(a_k^z)$
 $z_j^{\text{net}} = x_j$ $h^* = \sigma$ classification
 $\frac{\partial L}{\partial z_k^{\text{net}}} = \frac{\partial L}{\partial z_k^{\text{net}}} \frac{\partial z_k^{\text{net}}}{\partial z_k^{\text{net}}} = \frac{\partial L}{\partial z_k^{\text{net}}} h'(a_k^z) = \sum_n w_{kn} \frac{\partial L}{\partial z_n^{\text{net}}}$
 $\frac{\partial L}{\partial z_j^{\text{net}}} = h'(a_j^z) \sum_k w_{kj} z_k^{\text{net}}$
 $\frac{\partial L}{\partial w_{kj}} = \frac{\partial L}{\partial z_j^{\text{net}}} \frac{\partial z_j^{\text{net}}}{\partial w_{kj}} = \frac{\partial L}{\partial z_j^{\text{net}}} z_k^{\text{net}}$ backprop
 $\frac{\partial L}{\partial a_j^z} = \sum_k \frac{\partial L}{\partial z_k^{\text{net}}} \frac{\partial z_k^{\text{net}}}{\partial a_j^z} = \sum_k w_{kj} \frac{\partial L}{\partial z_k^{\text{net}}}$ recombine
 $p(y|x, w) = \prod_{n=1}^N p(y_n | x_n, w)^{1 - y_n} (1 - y_n)^{y_n}$ K separate bin. cls
 $E(w) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} + (1 - t_{nk}) \ln (1 - y_{nk})$ binary
 $E(w) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln (y_{nk})$
 $y_c(x) = \arg \max_i \left(\sum_j w_{ij} x_j \right)$
 $\sum_j \arg \max_i (z_j^{\text{net}})$
 $w^{\text{net}} = w^c + y^c DE$

Diagonal Approx
 $\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial y_j^{\text{net}}} z_j^{\text{net}}$ Hessian diagonal
 $\frac{\partial^2 E}{\partial w_{ij}^2} \sim h'(a_j^z) \sum_k w_{kj}^2 \frac{\partial^2 E}{\partial z_k^{\text{net}^2}} + h'(a_j^z) \sum_k w_{kj} \frac{\partial^2 E}{\partial z_k^{\text{net}} \partial a_j^z}$
 $O(w)$ by making diagonal Hessian approx
 $E(w) \sim E(w) + (w - \hat{w})^T \hat{b} + \frac{1}{2} (w - \hat{w})^T H (w - \hat{w})$
 $\hat{b} = \sqrt{E} w - \hat{w}$
 $\hat{w} E \sim \hat{b} + H (w - \hat{w})$
 $E(w) = E(w^*) + \frac{1}{2} \sum_i \lambda_i d_i^2$
 $\{u_i; u_i = w - w^*\}$
 also successive differentiating for $M_{1,1}^{-1} \text{Rank } M_{1,1}^{-1}$
 $H_{1,1} = \lambda_{1,1}$

Dropout
 → randomly 0 out weights + nodes remaining
 as steps so that leaves redundancy

Regularized
 → encourage weight sparsity
 → eigenvalue of Hessian $\ll \lambda \rightarrow \cup \lambda_i \rightarrow \frac{\lambda_i}{\lambda_i + \alpha}$

CNNs
 → translation invariant

$$S_{ij} = (I * K)_{ij} = \sum_{k,l} I_{ik} K_{lj}$$

↑
input

$$\sigma = \sum_{k,l} I_{ik} K_{lj} \quad (\text{cross-correlation})$$

Convolution \iff Circulant Matrix mul.

2D Conv \iff Doubly block circulant matrix

CNNs
 → typically quite sparse elements rather than full img - easier to train
 → param. sharing
 → equivariance: successive operations commute
 → f.g. equivariant: f.g. = gof
 → (translation) edge detection/scaling/translation commute (equivariant)

Pooling
 → max, avg
 → lazy form of compression
 → roughly invariant to small translations

Invariance
 → presence more useful than location
 → good for "summarizing"
 → or
 → global CNN robust param invariant
 → global pooling

Equivariance
 → able to identify in many orientations
 → encourage parameter sharing, makes respect symmetry, preserve "where"
 → permutational swap through in CNNs merge pairs
 → convolution

geometric CNNs have permutational equivariance
 but also translational/rotational/reflection equivariance

$$h_{ik}^e = \phi_k(h_{ik}^l, m_{ij}^e) \quad \begin{matrix} k = \text{feature} \\ ij = \text{node} \end{matrix}$$

$$m_{ik} = \sum_{j \neq i} g_{ij}$$

$$m_{ij}^e = \phi_e(h_{ik}^e, h_{jk}^e, a_{ij})$$

Equivariant Version

$$m_{ij}^e = \phi_e(h_{ik}^e, h_{jk}^e, \|x_i - x_j\|^2, a_{ij})$$

$$x_{ik}^e = x_{ik}^l + C \sum_{j \neq i} (x_{ik}^l - x_{jk}^l) \phi_x(m_{ij}^e)$$

Invariant: $f(Tx) = f(x)$
Equivariant: $f(Tx) = Tf(x)$

k-BH

→ k sets
 → all used to train + validate
 in (k-1), 1 node repeatedly
 → k=m (leave one out) high variance + costly to run
 → smaller = faster but more
 → large = lower bias