

Regression

Linear Regression: Model $P(w|X, \phi) = N_w[X^\top \phi, \sigma^2 I]$

Solution: Maximum Likelihood

$$\hat{\phi}, \hat{\sigma}^2 = \underset{\phi, \sigma^2}{\operatorname{argmax}} \left[-\frac{I \log(\sigma^2)}{2} - \frac{I \log(2\pi)}{2} + \frac{(w - X^\top \phi)^\top (w - X^\top \phi)}{2\sigma^2} \right]$$

$$\Rightarrow \hat{\phi} = (X^\top X)^{-1} X w$$

$$\sigma^2 = \frac{(w - X^\top \phi)^\top (w - X^\top \phi)}{I}$$

Bayesian Regression: Model $P(\phi|X, w) = \frac{P(w|X, \phi) P(\phi)}{P(w)}$

Closed-form formulation:

$$P(\phi|X, w) = N_\phi[\frac{1}{\sigma^2} A^{-1} X w, A^{-1}] \text{ with } A = \frac{1}{\sigma^2} X X^\top + \frac{1}{\sigma_p^2} I$$

$$\text{Reformulation: } A^{-1} = (\frac{1}{\sigma^2} X X^\top + \frac{1}{\sigma_p^2} I) = \sigma_p^2 I_D - \sigma_p^2 X (X^\top X + \frac{\sigma^2}{\sigma_p^2} I_1)^{-1} X^\top$$

$D < 1$ $D > 1$

$$\text{Likelihood: } P(w|X, \phi) = N_w[X^\top \phi, \sigma^2 I] = N_\phi[(\frac{1}{\sigma^2} X X^\top)^{-1} \frac{1}{\sigma^2} X w, (\frac{1}{\sigma^2} X X^\top)^{-1}]$$

$$\text{Prior: } P(\phi) = N_\phi[0, \sigma_p^2 I]$$

Inference: Integrate over the parameter space

$$\begin{aligned} P(w^*|x^*, X, w) &= \int P(w^*|x^*, \phi) P(\phi|X, w) d\phi \\ &= \int N_w[X^\top \phi, \sigma^2 I] N_\phi[\frac{1}{\sigma^2} A^{-1} X w, A^{-1}] d\phi \\ &= N_{w^*}[\frac{1}{\sigma^2} x^{*\top} A^{-1} X w, x^{*\top} A^{-1} x^* + \sigma^2] \end{aligned}$$

Fitting variance:

$$\begin{aligned} P(w|X, \sigma^2) &= \int P(w|X, \phi, \sigma^2) P(\phi) d\phi \\ &= \int N_w[X^\top \phi, \sigma^2 I] N_\phi[0, \sigma_p^2 I] d\phi \\ &= N_w[0, \sigma_p^2 X^\top X + \sigma^2 I] \end{aligned}$$

Non-linear regression: Model $P(w|z, \phi)$ with $z = \varphi(X)$

$$\hat{\phi} = (Z Z^\top)^{-1} Z w$$

$$\sigma^2 = \frac{(w - Z^\top \hat{\phi})^\top (w - Z^\top \hat{\phi})}{I}$$

Gaussian Process Regression:

Model:

$$P(w^* | x^*, X, w)$$

$$= \mathcal{N}_{w^*} \left[\frac{\sigma_p^2}{\sigma^2} k[x^*, X] w - \frac{\sigma_p^2}{\sigma^2} k[x^*, X] (k[X^*, X] + \frac{\sigma_p^2}{\sigma^2} I)^{-1} k[X, X] w, \right.$$

$$\left. \sigma_p^2 k[X^*, X^*] - \sigma_p^2 k[X^*, X] (k[X, X] + \frac{\sigma_p^2}{\sigma^2} I)^{-1} k[X, X^*] + \sigma^2 \right]$$

Sparse Linear Regression

Prior: $P(x) = \text{Stud}_x[\mu, \Sigma, v] = \frac{\Gamma(\frac{v+1}{2})}{\sqrt{v\pi\sigma^2}\Gamma(\frac{v}{2})} \left(1 + \frac{(x-\mu)^2}{v\sigma^2}\right)^{\frac{v-1}{2}}$

Approximation:

$$P(w|X, \sigma^2) = \max_H \left[\mathcal{N}_w[\mathbf{0}, X^T H^{-1} X + \sigma^2 I] \prod_{d=1}^D \text{Gam}_{h_d}[v/2, v/2] \right]$$

Learning: Learn H that maximizes above equation
Iteratively update posterior:

$$P(\phi|X, w) = \mathcal{N}_\phi \left[\frac{1}{\sigma^2} A^{-1} X w, A^{-1} \right] \text{ with } A = \frac{1}{\sigma^2} X X^T + H$$

$$\text{Estimate } H: h_d^{\text{new}} = \frac{1 - h_d \sum_{d' \neq d} v}{M_d^2 + v}$$

Dual Linear Regression: Model $P(w|X, \Theta) = \mathcal{N}_w[X^T X \psi, \sigma^2 I]$

Learning: Maximize Likelihood

$$\hat{\psi}, \hat{\sigma}^2 = \underset{\psi, \sigma}{\operatorname{argmax}} \left[\frac{I \log(2\pi)}{2} + \frac{I \log(\sigma^2)}{2} + \frac{(w - X^T X \psi)^T (w - X^T X \psi)}{2\sigma^2} \right]$$

$$\Rightarrow \hat{\psi} = (X X^T)^{-1} w$$

$$\hat{\sigma}^2 = \frac{(w - X^T X \psi)^T (w - X^T X \psi)}{I}$$

Dual Bayesian Regression:

Model: $P(\psi|X, w, \sigma^2) = \frac{P(w|X, \psi, \sigma^2) P(\psi)}{P(w|X, \sigma^2)} = \mathcal{N}_\psi \left[\frac{1}{\sigma^2} A^{-1} X^T X w, A^{-1} \right]$

$$\text{with: } A = \frac{1}{\sigma^2} X^T X X^T X + \frac{1}{\sigma_p^2} I$$

Relevance Vector Regression

Model $P(w|X, \Theta) = \mathcal{N}_w[X^T X \psi, \sigma^2 I]$

$$\text{with } P(\psi) = \prod_{i=1}^l \text{Stud}_i[\mathbf{0}, 1, v]$$

Learning: Same as sparse solution

Classification

Logistic Regression:

Model: $P(w|X, \phi) = \text{Bern} \left[\frac{1}{1 + \exp(-\phi^T x)} \right]$
 $= \left(\frac{1}{1 + \exp(-\phi^T x)} \right)^w \left(1 - \frac{1}{1 + \exp(\phi^T x)} \right)^{(1-w)}$

Log-likelihood:

$$L = \sum_{i=1}^I w_i \log \left[\frac{1}{1 + \exp(\phi^T x_i)} \right] + \sum_{i=1}^I (1-w_i) \log \left[\frac{\exp(\phi^T x_i)}{1 + \exp(\phi^T x_i)} \right]$$

Derivatives:

$$\frac{\partial L}{\partial \phi} = - \sum_{i=1}^I \left(\frac{1}{1 + \exp(\phi^T x_i)} - w_i \right) x_i = - \sum_{i=1}^I (\text{sig}[\phi^T x_i] - w_i) x_i$$

$$\frac{\partial^2 L}{\partial \phi^2} = - \sum_{i=1}^I \text{sig}[\phi^T x_i] (1 - \text{sig}[\phi^T x_i]) x_i x_i^T$$

Optimization: $\Theta^{+*} = \Theta^+ + \lambda \left(\frac{\partial^2 L}{\partial \phi^2} \right)^{-1} \frac{\partial L}{\partial \phi}$

Bayesian Logistic Regression:

Model: $P(\phi|X, w) = \frac{P(w|X, \phi) P(\phi)}{P(w|X)}$

Likelihood: $P(w|X, \phi) = \prod_{i=1}^I \left(\frac{1}{1 + \exp(\phi^T x_i)} \right)^{w_i} \left(\frac{\exp(\phi^T x_i)}{1 + \exp(\phi^T x_i)} \right)^{(1-w_i)}$

Prior: $P(\phi) = N_\phi[\phi_0, \phi_p^2 I]$

Learning: Find MAP solution using Newton's method

Log-posterior: $L = \sum_{i=1}^I \log[P(w|X, \phi)] + \log[P(\phi)]$

Derivatives: $\frac{\partial L}{\partial \phi} = - \sum_{i=1}^I (\text{sig}[\phi^T x_i] - w_i) x_i - \frac{\phi}{\phi_p^2}$

$$\frac{\partial^2 L}{\partial \phi^2} = - \sum_{i=1}^I \text{sig}[\phi^T x_i] (1 - \text{sig}[\phi^T x_i]) x_i x_i^T - \frac{1}{\phi_p^2}$$

Laplace Approximation: $P(\phi|X, w) \approx q(\phi) = N_\phi[\mu, \Sigma]$

Inference: $P(w^*|x^*, X, w) = \int P(w^*|x^*, \phi) P(\phi|X, w) d\phi$
 $\approx \frac{1}{1 + \exp[-\mu^T x^* / \sqrt{1 + \text{Tr}[(x^* \Sigma x^*) / 8}]}}$

Non-linear Logistic Regression

Model: $P(w|x, \phi) = \text{Bern}_w[\phi^T f(x)]$

Dual Logistic Regression!

Model: $P(w|X, \psi) = \text{Bern}_w[\text{sig}[X^\top \psi]]$

Derivatives: $\frac{\partial L}{\partial \psi} = -\sum_{i=1}^I (\text{sig}[\alpha_i] - w_i) X^\top x_i$

$\frac{\partial^2 L}{\partial \psi^2} = -\sum_{i=1}^I \text{sig}[\alpha_i](1 - \text{sig}[\alpha_i]) X^\top x_i x_i^\top X$

Relevance Vector Classification

Model: $P(w|X) = \int P(w|X, \psi) P(\psi) d\psi$

$$= \iint \prod_{i=1}^I \text{Bern}_{w_i}[\text{sig}[\psi^\top K[X_i, X_i]]] \mathcal{N}_{\psi_i}[\mathbf{0}, H^{-1}] \text{Gam}_{w_i}[\nu/2, \nu/2] dH d\psi$$

Prior: $P(\psi) = \prod_{i=1}^I \text{Stud}_{\psi_i}[\mathbf{0}, I, \nu]$

$$= \prod_{i=1}^I \int \mathcal{N}_{\psi_i}[\mathbf{0}, \frac{1}{h_i}] \text{Gam}_{w_i}[\frac{\nu}{2}, \frac{\nu}{2}] d\psi_i$$

$$= \int \mathcal{N}_{\psi_i}[\mathbf{0}, H^{-1}] \text{Gam}_{w_i}[\nu/2, \nu/2] dH$$

Random Forests

Splitting function: Axis-aligned: $h(v, \Theta) = \phi(v) \cdot \psi > z$
 with: $\phi(v) = (x_1, x_2)$
 $\psi = (0, 1) / (1, 0)$

Oriented line: $h(v, \Theta) = \phi(v) \cdot \psi > z$
 with: $\phi(v) = (x_1, x_2)$
 $\psi = (a, b) \in \mathbb{R}^2$

Conic section: $h(v, \Theta) = \phi(v)^T \psi \phi(v)$
 with: $\phi(v) = (x_1, x_2)$
 $\psi \in \mathbb{R}^{3 \times 3}$

Predictor model: Classification: Histogram over training samples

Regression: Constant: $y = \text{const.}$

Polynomial: $y = \sum_{i=0}^n w_i x^i$

Prob. linear: $y = L_1 x + L_2$

Density: $N(v; \mu_{L(v)}, \Lambda_{L(v)})$

Training: Classification: IG: $I = H(S_j) - \sum_{i \in LR} \frac{|S_j^i|}{|S_j|} H(S_j^i)$

Entropy: $H(S) = - \sum_{c \in C} p(c) \log(p(c))$

Regression: IG: $I = \sum_{(x,y) \in S_j} \log(\sigma_y(x)) - \sum_{i \in RL} \left(\sum_{(x,y) \in S_j^i} \log(\sigma_y(x)) \right)$

Entropy: $H(S) = \frac{1}{|S|} \sum_{x \in S} \int_y p(y|x) \log p(y|x) dy$
 $= \frac{1}{|S|} \sum_{x \in S} \frac{1}{2} \log((2\pi e)^2 \sigma_y^2(x))$

CART objective: $E = \sum_{(x,y) \in S_j} (y - \bar{y}_j)^2 - \sum_{i \in RL} \left(\sum_{(x,y) \in S_j^i} (y - \bar{y}_j)^2 \right)$

Density: IG: $I = H(S_j) - \sum_{i \in RL} \frac{|S_j^i|}{|S_j|} H(S_j^i)$

Entropy: $H(S) = \frac{1}{2} \log((2\pi e)^d |\Lambda(S)|)$

Ensemble model: $p(y|v) = \frac{1}{T} \sum_{t=1}^T p_t(y|v)$

Density: $p_t(v) = \frac{\pi_{L(v)}}{Z_t} N(v; \mu_{L(v)}, \Lambda_{L(v)})$

with $Z_t = \int \pi_{L(v)} N(v; \mu_{L(v)}, \Lambda_{L(v)}) dv$

Convergence theorem: $\mathbb{E}_{X, A(n), Z_t} \left[\left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\theta | X, A(n), Z_t] - \mathbb{E}[\theta | X] \right|^2 \right] \rightarrow 0$

for $k_n \rightarrow \infty, k_n/n \rightarrow 0, n \rightarrow \infty$

Maximum Margin Classification and Kernels

Classification Margin: $\gamma = \frac{\langle w, x_i \rangle + b}{\|w\|}$

Maximum Margin Classification

$\max_{\gamma \in \mathbb{R}, w \in \mathbb{R}^d} \gamma$ subject to: $y_i \langle w, x_i \rangle \geq \gamma$

$$\Leftrightarrow \min_{w \in \mathbb{R}^d} \|w\|^2 \text{ subject to: } y_i \langle w, x_i \rangle \geq 1$$

Soft-Margin Classification

$\min_{w \in \mathbb{R}^d, \xi_i \in \mathbb{R}^+} \|w\|^2 + C \sum_i \xi_i$ subject to: $y_i \langle w, x_i \rangle \geq 1 - \xi_i$

Hinge-Loss Formulation

$\min_{w \in \mathbb{R}^d} \|w\|^2 + C \sum_{i=1}^n (1 - y_i \langle w, x_i \rangle)_+$ with $[+]_+ = \max\{0, +\}$

Kernel Trick

$\min_{\alpha_i \in \mathbb{R}, \xi_i \in \mathbb{R}^+} \sum_{i,j} \alpha_i \alpha_j k[x_i, x_j] + C \sum_i \xi_i$ subject to: $y_i \sum_{j=1}^n \alpha_j k[x_i, x_j] \geq 1 - \xi_i$

Dual Formulation: Dualization with Lagrangian multipliers

$\min_{\alpha_i} \sum_{i,j} \alpha_i \alpha_j y_i y_j k[x_i, x_j] + C \sum_i \alpha_i$ subject to: $\sum_i y_i \alpha_i = 0$

Multiple Kernel Learning

$\min_{v_j \in \mathbb{R}, \xi_i \in \mathbb{R}^+, \beta_j \geq 0} \sum_j \frac{1}{\beta_j} \|v_j\|_{H_j}^2 + C \sum_i \xi_i$ subject to: $y_i \sum_j \langle v_j, \psi(x_i) \rangle_{H_j} \geq 1 - \xi_i$

Support Vector Regression

$\min_{w \in \mathbb{R}^d, \xi_i \in \mathbb{R}^+} \|w\|^2 + C \sum_{i=1}^n (\xi_i + \xi'_i)$ subject to: $y_i - \langle w, \psi(x_i) \rangle \leq \varepsilon + \xi_i$
 $\langle w, \psi(x_i) \rangle - y_i \leq \varepsilon + \xi'_i$

Dual Support Vector Regression

$\min_{\alpha_i \in \mathbb{R}, \xi_i \in \mathbb{R}^+} \sum_{i,j} \alpha_i \alpha_j k[x_i, x_j] + C \sum_i (\xi_i + \xi'_i)$ subject to: $y_i - \sum_j \alpha_j k[x_i, x_j] \leq \varepsilon + \xi_i$
 $\sum_j \alpha_j k[x_i, x_j] - y_i \leq \varepsilon + \xi'_i$

Multiclass SVM

Feature embedding: $\phi(x, y) = ([y=1] \phi(x), \dots, [y=k] \phi(x))$

$\min_{w \in \mathbb{R}^d, \xi_i \in \mathbb{R}^+} \frac{1}{2} \|w\|^2 + \frac{C}{N} \sum_{i=1}^N \xi_i$ subject to:
 $\forall i: \forall y \in \mathcal{Y} \langle w, \phi(x_i, y) \rangle - \langle w, \phi(x_i, y) \rangle \geq 1 - \xi_i$

Structured Support Vector Machine

Regularized Risk Minimization:

$$\min_{w \in \mathbb{R}^d} \lambda \|w\|^2 + \frac{C}{N} \sum_{i=1}^N \max_{y \in Y} [\Delta(y^i, y) + \langle w, \phi(x^i, y) \rangle - \langle w, \phi(x^i, y^i) \rangle]$$

Equivalent Formulations:

1. $\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{C}{N} \sum_{n=1}^N \xi_n$

subject to:

$$A_i: \max_{y \in Y} [\Delta(y^i, y) + \langle w, \phi(x^i, y) \rangle - \langle w, \phi(x^i, y^i) \rangle] \leq \xi_n$$

2. $\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{C}{N} \sum_{n=1}^N \xi_n$

subject to:

$$A_i: \forall y \in Y: \Delta(y^i, y) + \langle w, \phi(x^i, y) \rangle - \langle w, \phi(x^i, y^i) \rangle \leq \xi_n$$

3. $\min_{w, \xi} \|w\|^2 + \frac{C}{N} \sum_{n=1}^N \xi_n$

subject to:

$$A_i: \forall y: \langle w, \delta \phi(x^i, y^i, y) \rangle \geq \Delta(y^i, y) - \xi_n$$

4. $\min_{w, \xi} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n$

subject to:

$$A(\hat{y}^1, \dots, \hat{y}^N) \in Y \times \dots \times Y: \sum_{n=1}^N [\Delta(y^i, \hat{y}^n) + \langle w, \phi(x^i, \hat{y}^n) \rangle - \langle w, \phi(x^i, y^i) \rangle] \leq N \xi_n$$

Dual Structured Support Vector Machine

$$\max_{\alpha \in \mathbb{R}_{+}^{N \times |Y|}} \sum_{n=1, y \in Y}^N \alpha_{ny} \Delta(y^i, y) - \frac{1}{2} \sum_{\substack{y, \bar{y} \in Y \\ n, \bar{n} = 1, \dots, N}} \alpha_{ny} \alpha_{\bar{n}y} \langle \delta \phi(x^i, y^i, y), \delta \phi(x^{\bar{i}}, y^{\bar{i}}, y) \rangle$$

kernel: Function $K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$:

1. $K[x_1, x_2] = K[x_2, x_1]$

2. Kernel matrix $K_{ij} := (K[x_i, x_j])_{ij}$
is positive (semi-)definite: $\forall t \in \mathbb{R}^n: \sum_{i,j=1}^n t_i K_{ij} t_j \geq 0$
 \Leftrightarrow all eigenvalues $\lambda_i \geq 0$

Embedding function: $K[x, x'] = \langle \Phi(x), \Phi(x') \rangle_H$

kernels:

Polynomial: $K[x, x'] = (1 + \langle x, x' \rangle)^m$, $m > 0$

Gaussian: $K[x, x'] = \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right)$ $\sigma^2 \approx 0.6$

Set kernels

Averaging Kernels:

Counting: $K[D_i, D_j] = \sum_{i=1}^{n_i} \sum_{j=1}^{n_j} S_{d_i=d_j}$

Averaging: $K[D_i, D_j] = \frac{1}{n_i n_j} \sum_{i=j}^{n_i} \sum_{i=j}^{n_j} K[d_i, d_j]$

Harmonic: $K[D_i, D_j] = \left(\prod_{i=j}^{n_i} \prod_{i=j}^{n_j} K[d_i, d_j] \right)^{\frac{1}{n_i n_j}}$

Matching Kernel: $K[D_i, D_j] = \frac{1}{2} [k(D_i, D_j) + k(D_j, D_i)]$

with: $k(D_i, D_j) = \frac{1}{n_i} \sum_{i=1}^{n_i} \max_k K[d_i, d_j]$

Pyramid Match kernel:

$$K_A(\underbrace{\Psi(X), \Psi(Y)}_{\text{histogram pyramids}}) = \sum_{i=0}^L \frac{1}{2^i} \underbrace{\left(I(H_i(X), H_i(Y)) - I(H_{i-1}(X), H_{i-1}(Y)) \right)}_{\text{difficulty newly matched pairs}}$$

Histogram kernels

Unnormalized: **Linear kernel:** $K[h, h'] = \sum_j h_j h'_j$

Polynomial kernel: $K[h, h'] = (c + \sum_j h_j h'_j)^m$

Gaussian kernel: $K[h, h'] = \exp\left(-\frac{1}{\gamma} \sum_j \|h_j - h'_j\|^2\right)$

Normalized: **Histogram Intersection:** $K[h, h'] = \sum_j \min(h_j, h'_j)$

Bhattacharya: $K[h, h'] = \sum_j \sqrt{h_j h'_j}$

Symmetric kLD: $K[h, h'] = \exp\left(-\frac{1}{2} (KL(h|h') + KL(h'|h))\right)$

χ^2 kernel: $K[h, h'] = \exp\left(-\frac{1}{8} \chi^2(h, h')\right)$

with: $\chi^2(h, h') = \sum_j \frac{(h_j - h'_j)^2}{h_j + h'_j}$

Neural Networks

McCulloch-Pitts: $f(x, \Theta) = \sum_{d=1}^D \Theta_d x_d + \Theta_0$

Fully-connected Layer: $f(x, \Theta) = \sum_{n=0}^{N-1} \sum_{d=1}^D \Theta_d \phi(x_n) + \Theta_0$

Activation functions:

ReLU: $f(y) = \max(0, y)$

GELU: $f(y) = y \cdot P(Y \leq y) \approx 0.5y(1 + \tanh[\sqrt{2/\pi}(x + 0.044715y^2)])$

Sigmoid: $f(y) = \frac{1}{1 + \exp(-y)}$

Softmax: $f(y) = \frac{e^y}{\sum_y e^y}$

Backpropagation:

Error: $E = \frac{1}{2} \sum_{n=1}^N (\hat{y}_n - y_n)^2$

Weight change: $\Delta w_{ij} = \eta \cdot \delta_j \cdot z_i$

Output Neuron: $\Delta w_{ij} = -\eta \frac{\partial E}{\partial z_j} \cdot \frac{\partial z_j}{\partial \alpha_j} \cdot \frac{\partial \alpha_j}{\partial w_{ij}} = \eta(\hat{y}_j - y_j) \cdot f'(\alpha_j) \cdot z_j$

$$\delta_j = (\hat{y}_j - y_j) \cdot f'(\alpha_j)$$

Hidden Neuron: $\Delta w_{ij} = -\eta \left(\sum_{k=1}^K \delta_k \cdot \frac{\partial (\sum_h z_h w_{hk})}{\partial w_{jk}} \right) \frac{\partial z_i}{\partial \alpha_j} \frac{\partial \alpha_j}{\partial w_{ij}}$
 $= \eta \left(\sum_{k=1}^K \delta_k w_{jk} \right) f'(\alpha_j) z_j$

Bayesian NN: $R(\Theta) = \frac{1}{N} \sum_{n=0}^N L(y_n - f(x_n)) + \lambda \|\Theta\|^2$

Convolution: $(I * k)[x, y] = \sum_a \sum_b k[a, b] \cdot I[x-a, y-b]$

Segmentation Refinement

Energy function: $E(x) = \sum_i \varphi_u(x_i) + \sum_{i < j} \varphi_p(x_i, y_i)$

Unary: Inverse likelihood of segmentation

Pairwise: $\varphi_{ij}(x_i, x_j) = w_1 \exp\left(-\frac{\|p_i - p_j\|^2}{2\sigma_s^2} - \frac{\|I_i - I_j\|^2}{2\sigma_g^2}\right)$
 $+ w_2 \exp\left(-\frac{\|p_i - p_j\|^2}{2\sigma_\beta^2}\right)$

Variational Lower Bound:

$$\ln P(D) = \ln \int Q(\Theta) \frac{P(D, \Theta)}{Q(\Theta)} d\Theta \geq \int Q(\Theta) \ln P(D|\Theta) d\Theta + \int Q(\Theta) \ln \frac{1}{Q(\Theta)} d\Theta$$
$$= \ln P(D) - KL(Q(\Theta) \| P(\Theta | D))$$

Generative Models

Markov Random Fields

General potential function: $P(\omega) = \frac{1}{Z} \prod_{j=1}^J \phi_j(\omega_{C_j})$

MAP inference: $\hat{\omega}_{1\dots N} = \operatorname{argmin}_{\omega_{1\dots N}} \left[\sum_{n=1}^N U_n(\omega_n) + \sum_{(m,n) \in C} P_{mn}(\omega_m, \omega_n) \right]$

Binary Pairwise MRF: $P_\Theta(x) = \frac{1}{Z(\Theta)} \exp \left(\sum_{i,j \in E} x_i x_j \Theta_{ij} + \sum_{i \in V} x_i \Theta_i \right)$

Log-likelihood: $L(\Theta) = \frac{1}{N} \sum_{n=1}^N \log P_\Theta(x^{(n)})$

Derivative: $\frac{\partial L(\Theta)}{\partial \Theta_{ij}} = \frac{1}{N} \sum_n [x_i^{(n)} x_j^{(n)}] - \sum_x [x_i x_j P_\Theta(x)] = \mathbb{E}_{P_{\text{data}}}[x_i x_j] - \mathbb{E}_{P_\Theta}[x_i x_j]$

Restricted Boltzmann Machines

Energy function: $E(v, h; \Theta) = -\sum_{i,j} w_{ij} v_i h_j - \sum_i b_i v_i - \sum_j \alpha_j h_j$

Joint distribution: $P_\Theta(v, h) = \frac{1}{Z(\Theta)} \exp(-E(v, h; \Theta)) = \frac{1}{Z(\Theta)} \prod_{i,j} e^{w_{ij} v_i h_j} \prod_i e^{b_i v_i} \prod_j e^{\alpha_j h_j}$

with: $Z(\Theta) = \sum_{v,h} \exp(-E(v, h; \Theta))$

Generative distribution: $P_\Theta(v) = \sum_h P_\Theta(v, h)$

$= \frac{1}{Z(\Theta)} \prod_i e^{b_i v_i} \prod_j (1 + \exp(\alpha_j + \sum_i w_{ij} v_i))$

Model hidden variables: $P(h|v) = \prod_j P(h_j|v)$

with: $P(h_j=1|v) = \frac{1}{1 + \exp(-\sum_i w_{ij} v_i - \alpha_j)}$

Model visible variables: $P(v|h) = \prod_i P(v_i|h)$

with: $P(v_i=1|h) = \frac{1}{1 + \exp(-\sum_j w_{ij} h_j - b_i)}$

Maximum Likelihood: $L(\Theta) = \frac{1}{N} \sum_{n=1}^N \log P_\Theta(v^{(n)}) - \frac{\lambda}{N} \|w\|_F^2$

Derivative:

$$\begin{aligned} \frac{\partial L(\Theta)}{\partial w_{ij}} &= \frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial w_{ij}} \log \left(\sum_n \exp[v^{(n)\top} W h + \alpha^\top h + b^\top v^{(n)}] \right) - \frac{\partial}{\partial w_{ij}} \log Z(\Theta) - \frac{2\lambda}{N} w_{ij} \\ &= \mathbb{E}_{P_{\text{data}}}[v_i h_j] - \mathbb{E}_{P_\Theta}[v_i h_j] - \frac{2\lambda}{N} w_{ij} \end{aligned}$$

Weight update: $\Delta w_{ij} = \mathbb{E}_{P_{\text{data}}}[v_i h_j] - \mathbb{E}_P[v_i h_j]$

Gaussian-Bernoulli RBM

Energy function: $E(v, h; \Theta) = \sum_i \frac{(v_i - b_i)^2}{2\sigma_i^2} - \sum_{i,j} w_{ij} h_j \frac{v_i}{\sigma_i} - \sum_j \alpha_j h_j$

Model image: $P(v_i=x|h) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(\frac{(x_i - b_i - \Theta_i \sum_j w_{ij} h_j)^2}{2\sigma_i^2}\right)$

Model hidden variables: $P(h_j=1|v) = \frac{1}{1 + \exp(-\sum_i w_{ij} \frac{v_i}{\sigma_i} - \alpha_j)}$

Generative Adversarial Network

Discriminator Loss:

$$L(\theta) = \sum_j -\log [1 - \text{sig}[f[g(z_j, \theta), \phi]]] - \sum_i \log [\text{sig}[f[x_i, \phi]]]$$

$$\text{Generator Loss: } L(\theta) = \sum_j \log [1 - \text{sig}[f[g(z_j, \theta), \phi]]]$$

Jenson-Shannon divergence:

$$D_{JS}[P(x^*) \| P(x)]$$

$$= \frac{1}{2} \int P(x^*) \log \left[\frac{2P(x^*)}{P(x^*) + P(x)} \right] dx^* + \frac{1}{2} \int P(x) \log \left[\frac{2P(x)}{P(x) + P(x^*)} \right] dx$$

Wasserstein GAN

$$\text{Lipschitz continuous: } |D(x_2) - D(x_1)| \leq \|x_2 - x_1\|$$

$$\text{Lipschitz penalty: } R_{LP} = \mathbb{E}_{(x,h) \sim P_{x,h}} [\max(0, \|\nabla_{x,h} D(x,h)\|_2 - 1)^2]$$

$$\text{Wasserstein distance: } D_w[P(x) \| q(x)] = \min_P \left[\sum_{i,j} P_{ij} |i - j| \right]$$

$$\text{such that: } \sum_j P_{ij} = P(x=i)$$

$$\sum_i P_{ij} = q(x=j)$$

$$P_{ij} \geq 0$$

$$\text{Primal: } D_w[P(x), q(x)] = \min_{\Pi_{[T, \cdot]}} \left[\left\| \Pi_{[x_1, x_2]} \|x_1 - x_2\| dx_1 dx_2 \right\| \right]$$

$$\text{Dual: } D_w[P(x), q(x)] = \max_{f(x)} \left[\int P(x) f(x) dx - \int P(x^*) f(x^*) dx \right]$$

$$\text{subject to: } \|f(z_1) - f(z_2)\| \leq \beta \|z_1 - z_2\|$$

Diffusion Models

$$\text{Forward process: } q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1-\beta_t} x_{t-1}, \beta_t I)$$

$$q(x_{1:T} | x_0) = \prod_{t=1}^T q(x_t | x_{t-1})$$

$$\text{Closed-form: } q(x_T | x_0) = \mathcal{N}(x_T; \sqrt{\bar{\beta}_T} x_0, (1 - \bar{\beta}_T) I)$$

$$\text{with } \bar{\beta}_t = \prod_{s=1}^t (1 - \beta_s)$$

$$\text{Sampling: } x_T = \sqrt{\bar{\beta}_T} x_0 + \sqrt{1 - \bar{\beta}_T} \varepsilon$$

$$\text{Backward process: } p_\Theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\Theta(x_t, t), \sigma_\Theta^2 I)$$

$$p_\Theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\Theta(x_{t-1} | x_t)$$

Conditional:

$$p(z_{t-1} | z_t, x) = \mathcal{N}_{z_{t-1}} \left[\frac{(1 - \alpha_t)}{1 - \alpha_t} \sqrt{1 - \beta_t} z_t + \frac{\sqrt{\alpha_t} \beta_t}{1 - \alpha_t} x, \frac{\beta_t (1 - \alpha_t)}{1 - \alpha_t} I \right]$$

$$\text{Loss: } L(\phi) = \sum_{i=1}^T \sum_{t=1}^T \| \phi_t \left[\sqrt{\alpha_t} x_i + \sqrt{1 - \alpha_t} \cdot \varepsilon_i, \phi_t \right] - \varepsilon_i \|^2$$

Gaussian Processes

Linear Latent Variable Model: $y_i = w x_i + \varepsilon_i$ with $\varepsilon_i \sim N[0, \sigma^2]$

Probabilistic PCA: $p(Y|w) = \prod_{i=1}^n N(y_i | 0, w w^\top + \sigma^2 I)$

Maximum likelihood solution:

For: $\log p(Y|w) = -\frac{n}{2} \log |C| - \frac{1}{2} + r(C^{-1} Y Y) + \text{const.}$

with: $C = w w^\top + \sigma^2 I$

$$\Rightarrow w = U_q L R^\top, L = (\Lambda_q - \sigma^2 I)^{\frac{1}{2}}$$

with U_q : first q principal eigenvectors of $n^{-1} Y^\top Y$

Dual Probabilistic PCA: $p(Y|w) = \prod_{j=1}^P N(y_j | 0, X X^\top + \sigma^2 I)$

Maximum likelihood solution:

For: $\log p(Y|w) = -\frac{P}{2} \log |K| - \frac{1}{2} + r(K^{-1} Y Y) + \text{const.}$

with: $K = X X^\top + \sigma^2 I$

$$\Rightarrow w = U_q' L R^\top, L = (\Lambda_q - \sigma^2 I)^{\frac{1}{2}}$$

with U_q' : first q principal eigenvectors of $P^{-1} Y^\top Y$

Equivalence of solutions: $U_q = Y^\top U_q' \Lambda_q^{-\frac{1}{2}}$

Non-linear Latent Variable Models: $p(Y|X) = \prod_{j=1}^P N(y_j | 0, k)$

k : non-linear kernel

Character Animation: $\underset{x,q}{\operatorname{argmin}} \left(\frac{\|w(y - f(x))\|^2}{2\sigma^2(x)} + \frac{D}{2} \ln \sigma^2(x) + \frac{1}{2} \|x\|^2 \right)$
such that: $C(q) = 0$

Continuous Dimensionality Reduction:

Prior: $\sum_{i=1}^P \phi(s_i)$ with: $\phi(s_i, r) = |s_i|^r$ s_i : eigenvalues

Optimization: $\min_{Y, \Theta} p(Y|X, \Theta)$ s.t. $\forall i: s_i > 0, E(Y) - E(X) = 0$
with: $E(X) = \sum_i s_i^2$

Choice: $Q = \operatorname{argmax}_i \frac{s_i}{s_{i+1} + \varepsilon}$ with $\varepsilon \ll 1$ and $s_1 > s_2 > \dots > s_D$

Shared-Private Factorization

Optimization: $L = L_{\text{data}} + L_{\text{ortho}} + L_{\text{dim}} + L_{\text{energy}}$

$$L_{\text{ortho}} = 2 \sum_i (\|X^\top Z^{(i)}\|_F^2 + \sum_{j>i} \|Z^{(i)\top} Z^{(j)}\|_F^2)$$

Multilinear Models

Style-Content Separation: $y = \sum_{i,j} w_{ij} \alpha_i b_j + \varepsilon$

Multi-linear Analysis: $y = \sum_{i,j,k,\dots} w_{ijk\dots} \alpha_i b_j c_{k\dots} + \varepsilon$

Non-linear Basis Function: $y = \sum_{i,j} w_{ij} \alpha_i \phi_j(b) + \varepsilon$

Multi (non)-linear Models:

GP-LVM: $y = \sum_j w_j \phi_j(x) + \varepsilon = w^T \phi(x) + \varepsilon$

$$\text{with: } E[y, y'] = \phi(x)^T \phi(x') + \beta^{-1} S \\ = k[x, x'] + \beta^{-1} S$$

Multi-factor Gaussian Process: $y = \sum_{ijk\dots} w_{ijk\dots} \phi_i^{(1)} \phi_j^{(1)} \phi_k^{(1)} + \varepsilon$

$$\text{with: } E[y, y'] = \prod_i \phi^{(i)T} \phi^{(i)} + \beta^{-1} S \\ = \prod_i K_i[x^{(i)}, x'^{(i)}] + \beta^{-1} S$$

Continuous Character Control

Prior: $\ln P(X) = w_c \sum_{i,j} \ln k_{ij}^d$

Graph diffusion kernel: $k_{ij}^d = \exp(\beta H)$ with: $H = -T^{-1/2} L T^{-1/2}$

$T_{ii} = \sum_j w(x_i, x_j) \rightarrow \text{diagonal matrix and } w(x_i, x_j) = \|x_i - x_j\|^{-p}$

$$L_{ij} = \begin{cases} \sum_k w(x_i, x_k) & \text{if } i=j \\ -w(x_i, x_j) & \text{else} \end{cases}$$

Structured Learning

KL Divergence: $KL_{\text{cond}}(P||d)(x) = \sum_{y \in Y} d(y|x) \log \frac{d(y|x)}{P(y|x, w)}$

Expected KL: $KL_{\text{tot}}(P||d) = \mathbb{E}_{x \sim d(x)} [KL_{\text{cond}}(P||d)(x)]$
 $= \sum_{x \in X} \sum_{y \in Y} d(y|x) \log \frac{d(y|x)}{P(y|x, w)}$

Learning Techniques:

Maximum Likelihood:

$$\begin{aligned} w^* &= \underset{w \in \mathbb{R}^d}{\operatorname{argmax}} \sum_{n=1}^N P(y_1, \dots, y_N | x_1, \dots, x_N, w) \\ &= \underset{w \in \mathbb{R}^d}{\operatorname{argmax}} \prod_{n=1}^N P(y_n | x_n, w) \\ &= \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} - \sum_{n=1}^N \log P(y_n | x_n, w) \end{aligned}$$

Best Approximation:

$$\begin{aligned} w^* &= \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{x \in X} \sum_{y \in Y} d(x, y) \log \frac{d(y|x)}{P(y|x, w)} \\ &= \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} - \sum_{y \in Y} d(x, y) \log P(y|x, w) \\ &= \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} - \sum_{n=1}^N \log P(y_n | x_n, w) \end{aligned}$$

MAP Estimation: $w^* = \underset{w \in \mathbb{R}^d}{\operatorname{argmax}} P(w|D) = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} [-\log P(w|D)]$
 $= \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} [-\log P(w) - \sum_{n=1}^N \log P(y_n | x_n, w) + \log P(y_n | x_n)]$
 $= \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} [-\log P(w) - \sum_{n=1}^N \log P(y_n | x_n, w)]$

Prior: $P(w) = \text{const.} \Rightarrow w^* = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} [-\sum_{n=1}^N \log P(y_n | x_n, w) + \text{const.}]$
 $P(w) = \text{const.} \cdot e^{-\frac{1}{2\sigma^2} \|w\|^2}$
 $\Rightarrow w^* = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} [\frac{1}{2\sigma^2} \|w\|^2 - \sum_{n=1}^N \log P(y_n | x_n, w) + \text{const.}]$

CRF Learning

Model: $P(y|x, w) = \frac{1}{Z(x, w)} \exp(\langle w, \phi(x, y) \rangle)$

Likelihood: $L(w) = \frac{1}{2\sigma^2} \|w\|^2 - \sum_{n=1}^N [\langle w, \phi(x_n, y_n) \rangle - \log \sum_{y \in Y} e^{\langle w, \phi(x_n, y) \rangle}]$

Gradient: $\nabla_w L(w) = \frac{1}{\sigma^2} w - \sum_{n=1}^N [\phi(x_n, y_n) - \mathbb{E}_{y \sim P(y|x_n, w)} [\phi(x_n, y)]]$

Hessian: $\Delta L(w) = \frac{1}{\sigma^2} I_{D \times D} + \sum_{n=1}^N (\mathbb{E}_{y \sim P(y|x_n, w)} [\phi(x_n, y)]) (\mathbb{E}_{y \sim P(y|x_n, w)} [\phi(x_n, y)])^T$