# Exponential Family

FAMILY OF PROB. DISTRIBUTIONS WITH VERY COOL PROPERTIES. GAUSSIAN, BERNOULLI, STUDENT'S T

- ONLY FAMILY WITH FINITE SUFFICIENT STATISTICS → MAKES LOSSLESS COMPRESSION FEASIBLE → REQUIRES SUPPORT NOT BE
  UNDER REGULARITY                                                                              DEPENDANT FROM PARAMETERS
- ONLY FAMILY WITH CONJUGATE PRIORS
- LEAST SET OF ASSUMPTIONS UNDER CONSTRAINS
- USED FOR GLM, VARIATIONAL INFERENCE

## DEFINITION

A PDF/PMF IS IN EXP FAMILY IF

$$P(x|\theta) = \frac{1}{Z(\theta)} h(x) \exp[\theta^T \phi(x)] = h(x) \exp[\theta^T \phi(x) - A(\theta)]$$

$$Z(\theta) = \int h(x) \exp[\theta^T \phi(x)] dx \qquad A(\theta) = \log(Z(\theta))$$

$\phi(x) =$ VECTOR OF SUFFICIENT STATISTICS

$Z(\theta) =$ PARTITION FUNCTION

$h(x) =$ SCALING, OFTEN $=1$   • IF $\phi(x) = x \to$ NATURAL EXPONENTIAL FAMILY

$A(\theta) =$ LOG PARTITION FUNCTION / CUMULANT FUNCTION GENERATING   $\eta(\theta):$ PARAMETERS $\longrightarrow$ CANONICAL PARAMETERS

• IF $\dim(\theta) < \dim(\eta(\theta))$ CURVED EXPONENTIAL FAMILY

• IF $\eta(\theta) = \theta \to$ CANONICAL FORM

## BERNOULLI

$$BER(x|\mu) = \mu^x (1-\mu)^{1-x} = \exp[x \log(\mu) + (1-x)\log(1-\mu)] = \exp[\phi(x)^T \theta], \quad \phi(x) = [I(x=0), I(x=1)], \quad \theta = [\log(\mu), \log(1-\mu)]$$

MINIMAL REPRESENTATION $\longrightarrow$ UNIQUE $\theta$ ASSOCIATED          IS OVERCOMPLETE, LINEAR DEPENDENCE

$$BER(x|\mu) = (1-\mu) \exp\left[x \cdot \log\left(\frac{\mu}{1-\mu}\right)\right] \qquad \phi(x) = x, \quad \theta = \log\left(\frac{\mu}{1-\mu}\right) \text{ LOG ODDS RATIO} \quad . \quad Z = \frac{1}{1-\mu} \left.\right\} \text{ CANONICAL FORM}$$

$$\mu = \text{SIGM}(\theta) = \frac{1}{1+e^{-\theta}}$$

## UNIVARIATE GAUSSIAN

$$N(x|\mu, \sigma^2) = \frac{1}{Z(\theta)} \exp(\theta^T \phi(x)) \qquad \theta = \begin{pmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{pmatrix} \quad \phi(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix} \quad Z(\mu, \sigma^2) = \sqrt{2\pi\sigma} \exp\left[\frac{\mu^2}{2\sigma^2}\right]$$

## CUMULANTS

DERIVATIVES OF $A(\theta) \to$ CUMULANTS

$$\frac{dA}{d\theta} = \frac{d}{d\theta}\left(\log \int \exp(\theta \phi(x)) h(x) dx\right) = \int \phi(x) p(x) dx = E[\phi(x)]$$

$$\frac{d^2A}{d\theta} = \ldots = E[\phi^2(x)] - E[\phi(x)]^2 = VAR(x)$$

IF MULTIV: $\nabla^2 A(\theta) = COV(\phi(x))$
$\to A(\theta)$ IS CONVEX

FOR BERNOULLI:  $A(\theta) = \log(1+e^\theta)$,  $\frac{dA}{d\theta} = \frac{e^\theta}{1+e^\theta} = \frac{1}{1+e^{-\theta}} = \text{SIGM}(\theta) = \mu$,  $\frac{d^2A}{d\theta} = (1-\mu)\mu$

## MLE

LIKELIHOOD: $P(D|\theta) = \left[\prod_1^N h(x_i)\right] g(\theta)^N \exp\left(\eta(\theta)^T \cdot \left|\sum_i^N \phi(x_i)\right|\right)$, $\phi_D = \left[\sum_i \phi_i(x_i) \dots \sum_N \phi_N(x_i)\right]$

LOG-LIKELIHOOD: $\log(P(D|\theta)) = \theta^T \phi(D) - NA(\theta)$. • $-A(\theta)$ is CONCAVE, $\theta^T \phi(x)$ LINEAR $\longrightarrow$ LOG LIKELIHOOD IS CONCAVE, UNIQUE GLOBAL M

$$\nabla_\theta \log P(D|\theta) = \phi(D) - N E[\phi(x)]$$

$$E[\phi(x)] = \frac{1}{N} \sum \phi(x_i)$$ • EMPIRICAL AVERAGE OF SUFF. STATS. MUST EQUAL THEORETICAL XPECTED SUFF. STATS,

MEHODS OF MOMENTS (2)

IE FOR BERNOULLI: $E[\phi(x)] = P(x=1) = \hat{\mu} = \frac{1}{N} \sum_i^N \mathbb{I}(x=1)$

## BAYESIAN FORMULATIONS

FOR CONJUGATE PRIOR TO MAKE SENSE $\longrightarrow$ LIKELIHOOD MUST HAVE FINITE STATISTICS $\longrightarrow$ ONLY EXP FAMILY

LIKELIHOOD: $P(D|\theta) \propto g(\theta)^N \exp(\eta(\theta)^T S_N)$, $S_N = \sum S$. $\left|\begin{array}{c} P(D|\eta) \\ \propto \end{array} \exp\left(N\eta^T \bar{S} - NA(\theta)\right)\right.$

PRIOR: $P(\theta|v_0, \tau_0) \propto g(\theta)^{v_0} \exp(\eta(\theta)^T \tau_0)$ | $P(\eta|v_0, \tau_0) = \exp\left(v_0 \eta_0^T \bar{\tau}_0 - v_0 A(\eta)\right)$

POSTERIOR: $P(\theta|D) = P(\theta|v_N, \tau_N) = P(\theta|v_0+N, \tau_0+S_N)$ | $P(\eta|D) = P\left(\eta\big|v_0+N, \frac{v_0\bar{\tau}_0 + N\bar{S}}{v_0+N}\right)$ CONVEX COMBINATION OF
$\underset{\text{JUST UPDATES HYPER.PARAMS BY}}{\underset{\text{ADDITION}}{}}$ PRIOR MEAN HYPERPARAMS AND AVG OF SUFF. STATS

POSTERIOR PREDICTIVE: $P(D'|D) = \int P(D'|\theta) P(\theta|D) d\theta = \left[\prod_1^N h(\bar{x}_i)\right] \frac{z(\bar{\tau}_0 + \bar{S}(D) + \bar{S}(D'))}{z(\bar{\tau}_0 + \bar{S}(D))}$

## MAXIMUM ENTROPY DERIVATION

MAXIMUM ENTROPY PRINCIPLE: WE SHOULD PICK DISTRIBUTION WITH MAX ENTROPY (CLOSEST TO UNIFORM), IMPOSING THAT MOMENTS MATCH WITH EMPIRICAL
MOMENTS OF FUNCTION
$\langle$ SOME CALCULATIONS LATER $\rangle$ $P(x) = \frac{1}{z} \exp\left(-\sum \lambda_n f_n(x)\right)$ HAS FORM OF EXP FAMILY, GIBBS DISTRIBUTION

# Generalized Linear Models

ANY MODEL WHERE OUTPUT DENSITY IS IN EXPONENTIAL FAMILY AND WHERE MEAN PARAMS ARE LINEAR COMBINATION OF INPUTS, POSSIBLY THROUGH A NONLINEARITY. LINEAR REGRESSION AND LOGISTIC REGRESSION ARE GLM.

## BASICS

$\sigma^2$ IS DISPERSION PARAMETER, $\theta$ IS CANONICAL PARAMETER, $\mu$ IS MEAN PARAMETER. TO GO FROM $\theta$ TO $\mu$ WE USE $\psi \to \theta$, $\theta = \psi(\mu)$, NORMAL IS INVERTIBLE

MEAN FUNCTION: $\eta_i = w^T x_i \longrightarrow \mu_i = g^{-1}(\eta_i) = g^{-1}(w^T x)$

$\mu = \psi^{-1}(\theta) = A'(\theta)$

LINK FUNCTION: $g()$, CAN PICK ANY AS LONG AS IT'S INVERTIBLE. IN LOGISTIC $g^{-1} = $ SIGM

CANONICAL LINK FUNCTION: IF WE PICK $g = \psi$. ONE IN BERNOULLI $g(\mu) = \log\left(\dfrac{\eta}{1-\eta}\right)$ LOGIT FCN, $\mu = $ SIGM$(\eta)$

ALSO, VERY VERY GENERALLY, $E[y|x, w, \sigma^2] = \mu_i = A'(\theta)$

$\text{var}[y|x, w, \sigma^2] = \sigma_i^2 = A''(\theta)\sigma^2$

- NORMAL $\longrightarrow$ IDENTITY $\longrightarrow \theta = \mu$
- POISSON $\longrightarrow$ LOG $\longrightarrow \theta = \log(\mu)$

## ML AND MAP

ALL GLM CAN BE FIT W/ SAME PROCEDURE AS LOGISTIC REGRESSION $\ell(w) = \log p(D|w) = \dfrac{1}{\sigma^2}\sum \ell_i$, $\ell_i = \theta_i y_i - A(\theta_i)$

GRADIENT WITH CHAIN RULE $\dfrac{d\ell_i}{dw_j} = \dfrac{d\ell}{d\theta} \cdot \dfrac{d\theta}{d\mu_i} \cdot \dfrac{d\mu_i}{d\eta} \cdot \dfrac{d\eta}{dw_j} = (\ldots)$ IF CANONICAL LINK $\theta_i = \eta_i$, $\nabla_w \ell(w) = \dfrac{1}{\sigma^2}\left[\sum_i^N (y_i - \mu_i)x_i\right]$

- THEN WE CAN USE GRADIENT DESCENT ALGORITHM

SUM OF INPUT VECTOR WEIGHED BY ERRORS.

- OR BETTER, A 2ND ORDER METHOD $H = -\dfrac{1}{\sigma^2}X^T S X$, S IS DIAGONAL WEIGHING MATRIX (CANONICAL LINK)
  - IF NON-CANONICAL LINK, USE EXPECTED HESSIAN, OR FISHER INFORMATION MATRIX, HAS SAME FORM OF H UNDER CANONICAL LINK
- IF MAP: INTRODUCE A GAUSSIAN PRIOR, LIKE $L^2$ REGULARIZATION IN LOGISTIC REGRESSION

BAYESIAN INFERENCE: MCMC, VARIATIONAL INFERENCE, OR GAUSSIAN APPROXIMATIONS

# PROBIT REGRESSION

$g^{-1}(\eta) = \Phi(\eta)$, ERF, CDF OF GAUSSIAN

- FIND GRADIENT, HESSIAN, STILL CLOSED FORM
- PLUG IN GRADIENT-BASED OPTIMIZER
- CAN BE INTERPRETED AS RANDOM UTILITY MODEL (RUM)

- IS SUITED TO ORDINAL REGRESSION, WHERE RESPONSE IS DISCRETE-VALUED WITH ORDER. MULTIPLE THRESHOLDS.
- MULTINOMIAL PROBIT: UNORDERED CATEGORICAL VALUES; MODELS C CORRELATED BINARY OUTCOMES

# GLM 4 MULTI-TASK LEARNING

FIT MANY RELATED CLASSIFICATION MODELS. BETTER PERFORMANCE IF WE FIT ALL PARAMS AT SAME TIME. ALSO: **TRANSFER LEARNING** VIA **HIERARCHICAL BAYESIAN METHODS**                                                                    **LEARNING TO LEARN**

TIPS COLLABORATIVE FILTERING , MANY GROUPS, MANY FEATURES. MAJORITY OF GROUPS HAVE LITTLE DATA, LONG TAILS.
CAN'T FIT SAME MODEL FOR ALL GROUPS , BUT CAN'T FIT EACH ONE SEPARATELY RELIABLY . WE **ENCOURAGE PARAMS TO BE SIMILAR**

- $E[y_{ij}|x_{ij}] = g(x_{ij}^T, \beta_j)$ , $\beta_j \sim N(\beta_*, \sigma_j^2 I)$ , $\beta_* \sim N(\mu, \sigma_*^2 I)$

  - CAN MAP WITH STD GRADIENT METHODS

GROUPS WITH SMALLER SAMPLE SIZE BORROW PREDICTIVE STRENGTH
FROM LARGER ONES VIA COMMON PARENT. $\beta_*$

  - $\sigma_j^2$ CONTROLS HOW MUCH DEVIANCE FROM C.P.
  - $\sigma_*^2$ OVERALL PRIOR STRENGTH

- **EXAMPLE** : PERSONALIZED MAIL FILTERING ⟶ $\beta_*$ FROM EVERYONE'S MAIL, $\beta_j$ FROM SINGLE USER MAIL
- OTHER PRIORS (THAN GAUSSIAN) POSSIBLE , ⟶ SPARSITY INDUCING PRIOR ON $\beta_j$ FOR MULTI-TASK FEATURE SELECTION (CONJOINT ANALYSIS)
- **NEGATIVE TRANSFER**: MULTITASK LEARNING FUCKS UP (WORSE) WHEN PARAMS ARE QUALITATIVELY DIFFERENT, WRONG INDUCTIVE BIAS ON PRIOR

# GENERALIZED LINEAR MIXED MODELS

INFORMATION AT BOTH GROUP LEVEL AND ITEM LEVEL, STILL MULTITASK , LIKE ANOVA A BIT.

$$E[y_{ij}|x_{ij}, x_j] = g(\phi_1(x_{ij})^T \beta_j + \phi_2(x_{ij})^T \alpha)$$    $\beta_j$ RANDOM FX    $\alpha$ FIXED FX    • IF $P(y|x)$ GLM ⟶ GLMM

VARY RANDOMLY ON GROUPS

- CAN BE DIFFICULT TO FIT ⟶ $P(y_{ij}|\theta)$ MAY NOT BE CONJUGATE TO $P(\theta)$, TWO LEVELS OF UNKNOWNS $\theta$ AND $\eta = (\mu, \sigma)$ FOR PRIOR.
- FULL BAYESIAN INFERENCE: MCMC, VARIATIONAL BAYES OR EMPIRICAL BAYES , OR EXPECTATION - MINIMIZATION

# RANKING

LEARNING TO RANK, USER PROBLEM. QUERY, DOCUMENTS, RELEVANCE, ETC...

- PROBABILISTIC LANGUAGE MODEL, BAG OF WORDS $\;\; sim(q,D) \triangleq P(q|D) = \prod P(q_i | D)$ . $q$ IS WORD, $D$ IS DOCUMENT
- HAS TO BE SMOOTHED $\;\; P(t|d) = (1-\lambda) \frac{TF(t,d)}{LEN(d)} + \lambda P(t | BACKGROUND)$ \;\; TF IS TERM FREQUENCY

## — POINTWISE

FOR EACH QUERY-DOCUMENT PAIR DEFINE FEATURE VECTOR $x(Q,D)$, LABELS (Y/N, OR ORDERED CATEGORICAL)

THEN WE $P(y=1 | x(Q,D))$ OR $P(y=R | x(Q,D))$ • SIMPLE BUT NO CONSIDER LOCATION OF DOCS IN RESULT LIST.

         • ERRORS AT TOP = ERRORS AT BEGINNING.

## — PAIRWISE

- $P(y_{iu} | x(Q, d_i), x(Q, d_u))$ \;\; $y=1 \rightarrow REL(d_i, q) > REL(d_u, q)$ ELSE $0$ \;\; BINARY CLASSIFIER
- $P(y_{iu}=1 | x_i, x_u) = SIGM(f(x_i) - f(x_u))$ ; F IS OFTEN LINEAR SCORING FUNCTION $f = w^T x$ \;\; **RANKNET**
- MLE OF $w$ VIA MAX LL OR MIN CROSS ENTROPY LOSS , OPTIMIZE W/ GRADIENT DESCENT

## — LISTWISE

- FULL CONTEXT FOR RELEVANCY
- DEFINE TOTAL ORDER ON LIST WITH PERMUTATION OF INDICES $\pi$ . PLACKETT-LUCE DISTR: $P(\pi | s) = \prod_1^M \frac{s_j}{\sum_u^M s_u}$ \;\; $s_j = (\pi^{-1}(j))$ SCORE OF DOC AT $j$ POSITION
- $\pi = (A,B,C) \rightarrow P(\pi) = P(A=1) P(B=2 | A=1) P(C=3 | B=2, A=1)$ • INCORPORATE FEATURES VIA $s(d) = f(x(Q,D))$ \;\; LINEAR $f = w^T x$
- MINIMIZE CROSS-ENTROPY $-\sum_i \sum_\pi P(\pi | y_i) \log P(\pi | s_i)$ , INTRACTABLE $\rightarrow$ CONSIDER ONLY <u>TOP K POSITIONS</u> \;\; **LISTNET**
- $K=1$ CROSSENTROPY TAKES $O(m)$ TIME
- IF ONLY 1 DOCUMENT IS RELEVANT $\rightarrow$ CAN USE MULTINOMIAL LOGISTIC / SOFTMAX $\;\; P(y=c | x) = \frac{EXP(s_c)}{\sum_1^m EXP(s_{c'})}$ \;\; USED IN COLLAB FILTERING

# ASSORTED LOSS FUNCTIONS FOR RANKING

- MEAN RECIPROCAL RANK (MRR) \;\; QUERY Q \;\; RANK OF FIRST $R(q)$ \;\; $MRR = 1/R(q)$
- MEAN AVG PRECISION (MAP) \;\; PRECISION $= P@K(\pi) =$ NUM RELEVANTS IN TOP K POS OF $\pi$ / $K$ , $AP(\pi) = \frac{\sum_K P@K(\pi)}{NUM\ RELEVANTS}$

  $MAP = \sum_q AP(\pi) / N_q$

- NORMALIZED DISCOUNTED CUMULATIVE GAIN (NDCG) \;\; RELEVANCE LABELS, MULTIPLE LEVELS

  $DCG@K(R) = R_1 + \sum_2^K \frac{R_i}{lg_2 i}$ , $R_i$ RELEVANCE • DCG VARIES WITH LENGTH OF LIST $\rightarrow$ NORMALIZED WITH OPTIMAL ORDERING $IDCG@K = \max_\pi DCG@K$

         • NDCG = ~~DCG/DCG~~ DCG/IDCG

- RANK CORRELATION BETWEEN RANKED LIST $\pi$ AND RELEVANCE JUDGEMENT ~~$\pi^*$~~ VIA I.E. **WEIGHTED KENDALL $\tau$ STATISTIC**
- **WARP LOSS** WEIGHTED APPROXIMATE PAIRWISE. BETTER THAN PRECISION @K TRANSFORMS INTEGER RANK TO REAL PENALTY

- LOSSES USED <u>BAYESIANLY</u> • FIT MODEL W/ POSTERIOR INFERENCE $\rightarrow$ THEN CHOSE ACTIONS TO MINIMIZE EXPECTED FUTURE LOSS. SAMPLE FROM POSTERIOR THEN AVG OVER $\beta$s FOR DIFFERENT THRESHOLDS

- <u>FREQUENTISTLY</u> MINIMIZE EMPIRICAL LOSS ON TRAINING SET, BUT NOT DIFFERENTIABLE $\rightarrow$ USE GRADIENT FREE OPTIMIZATION OR SURROGATE LOSSES, I.E CROSSENTROPY