

设定触发或过滤阈值的一个方法

臧石磊

slz2008@gmail.com

南京大学

粒子物理计算软件与技术研讨会

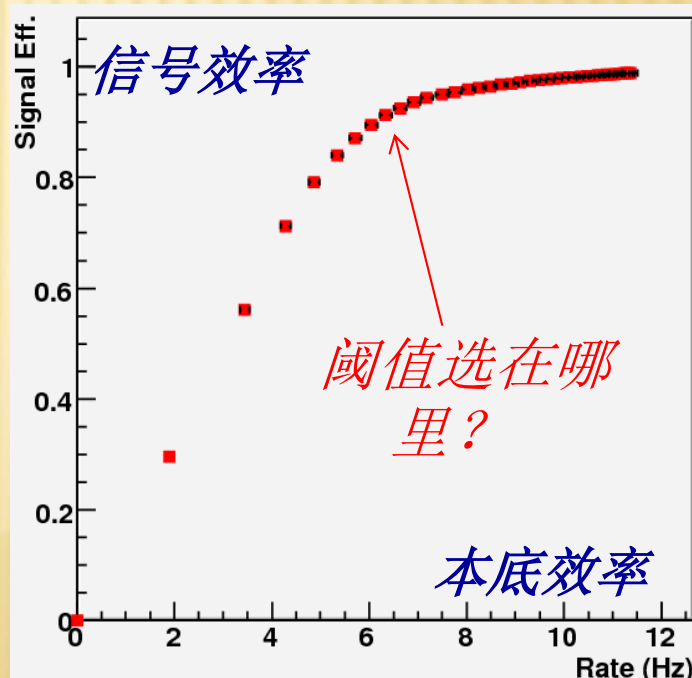
山东·威海

2013年7月4日

优化触发/过滤阈值

- **问题:** 在触发(trigger), 过滤(skim or filter), 粒子鉴别 (PID), 物理分析中, 如何设定阈值或cut值?
- **本报告:** 几年前的结果, 基于CMS实验的模拟数据, 优化超对称物理分析中的光子触发阈值(HLT)和分析中的光子选择条件。
- **本报告目的:** 和大家交流, 希望大家尝试和检验这个方法。

□ 信号效率 vs. 本底效率 →
怎样用一个客观的方法优化阈值?



问题来源

- 在触发、过滤、粒子鉴别、多变量分析和物理分析中，用一些变量来压制本底同时保留信号。**原因一**：资源限制（电子学、存储空间、计算能力CPU等）；**原因二**：通过压低本底得到清晰信号。
- **对于原因一**，一般通过研究信号效率vs本底效率，用眼观察设定一个合理的阈值。
- **对于原因二**，在传统物理分析中，为了得到清晰信号，会采用：
 - 选择条件或cut条件的设置，用和物理量相关的统计量来优化：*分支比的误差最小；结果的显著性最大；90%置信上限最好；等等。即：*

$$\text{Significance} \equiv \sqrt{2\ln Q} \approx 2(\sqrt{N_S + N_B} - \sqrt{N_B}) \approx N_S / \sqrt{N_B} \approx N_S / \sqrt{N_S + N_B}$$

- 在质量测量中（如top质量测量），*最小化测量质量的误差。*
- *→ 缺点：测量结果被人为地调试得更好。测量结果不全是盲(blind)分析。结果有人为倾向性。*

信息理论(I)

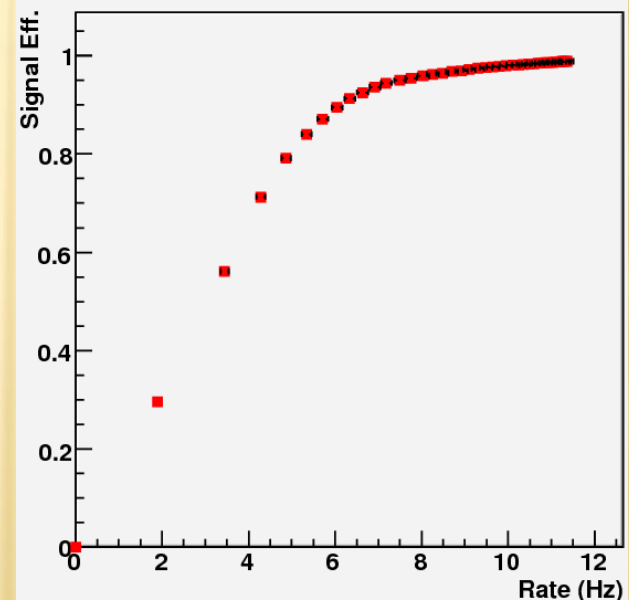
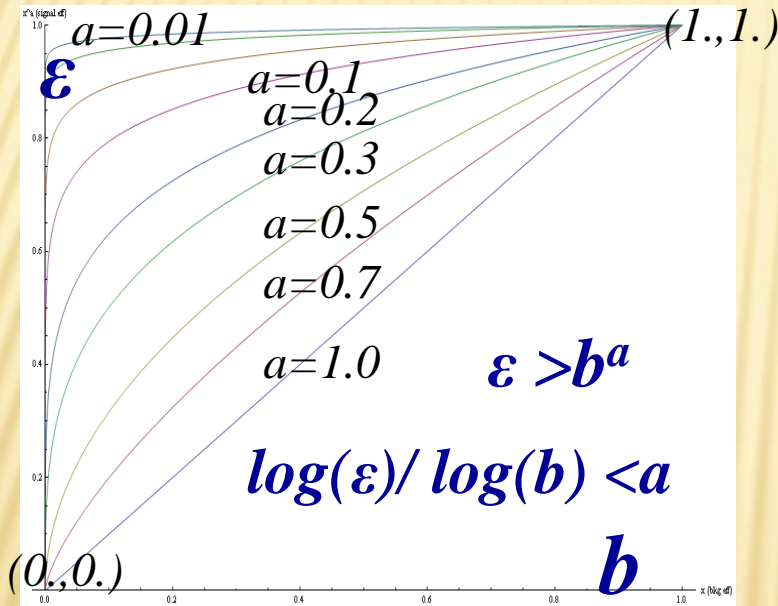
- 1948年香农(C. Shannon)提出信息理论
- N 个数（事例数）：信息量: $\log_2 N$.
- 信号和本底的信息量: $\log(N_S)$, $\log(N_B)$
- 一个cut带来的信号效率 ϵ 和本底效率 b
- Cut之后的信息量: $\log(N_S \epsilon)$, $\log(N_B b)$
- 信息量的减少: $-\log(\epsilon)$, $-\log(b)$
- ✓ 信号和本底信息量减少的比值: $\log(\epsilon)/\log(b)$
- 假设两个cut条件减少的信息量相同，即
 $-\log(\epsilon)-\log(b) = -\log(\epsilon')-\log(b')$,
如果 $\log(\epsilon)/\log(b) < \log(\epsilon')/\log(b')$ \rightarrow 那么cut 比cut'好
- $\log(\epsilon)/\log(b)$ 越小越好
- $\log(\epsilon)/\log(b) < a \rightarrow \epsilon > b^a$ ($0 < \epsilon, b, a \leq 1$).
 \rightarrow 统计量 $\log(\epsilon)/\log(b)$ 做为优化标准

信息理论(II)

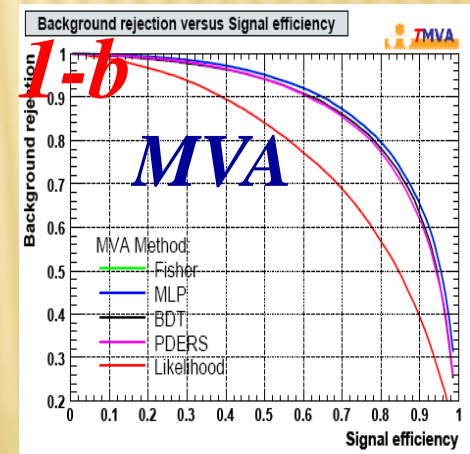
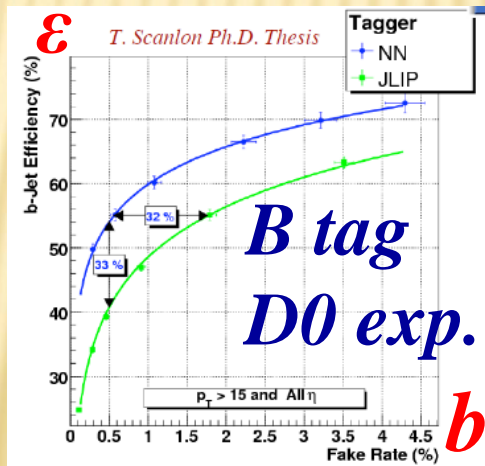
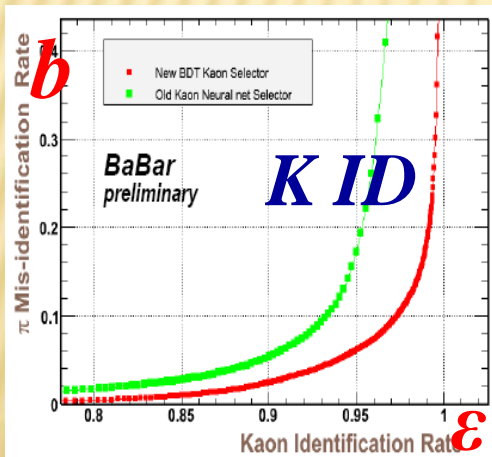
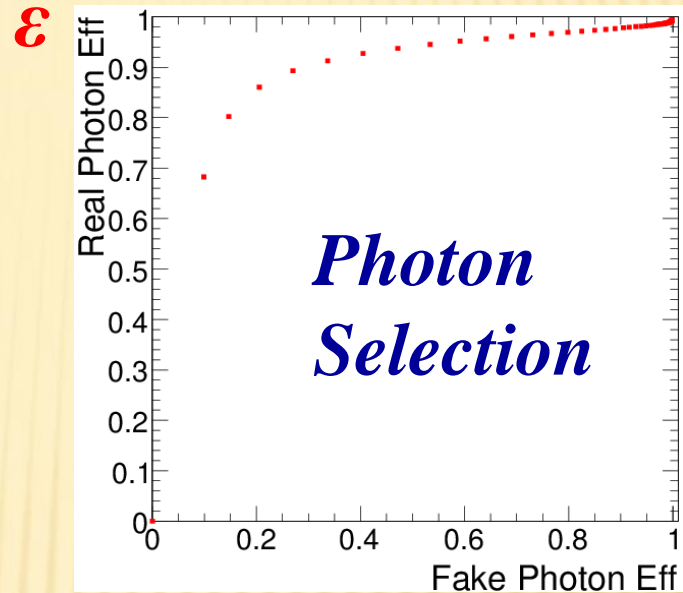
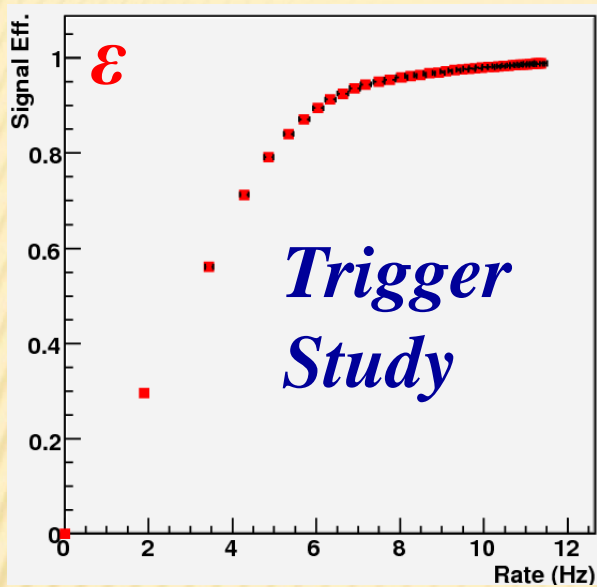
- N 是信使数量，物理结果是N个信使携带的消息的意义。
- 对于**分支比**，信使数量就是信使携带的消息的意义。
- 对于**宽度、质量、极化**等等物理量，消息的意义不仅仅是信使数量，还有信使所携带的其他信息（例如：动量、能量、角度、顶点等等）。
- **优点：盲分析！**
 $\log(\epsilon)/\log(b)$ 只涉及信息量，不涉及消息的意义。
- **注意：**在其他cut条件都使用的情况下，用 $\log(\epsilon)/\log(b)$ 优化一个cut条件。
 - 信息量的定义满足如下属性： $\log_2(xy) = \log_2(x) + \log_2(y)$ ，并且是满足这个属性的唯一解。
 - 这个属性对 $\log(\epsilon)/\log(b)$ 方法非常重要，假设有两个cut条件， $\log(\epsilon_1 \epsilon_2)/\log(b_1 b_2) = (\log(\epsilon_1) + \log(\epsilon_2))/(\log(b_1) + \log(b_2)) \rightarrow$ 两个cut分开优化和合起来优化产生一致的结果。

$\log(\epsilon)/\log(b)$ 方法

- 信号和本底信息量减少的比值 $\log(\epsilon)/\log(b)$ 做为优化标准
- $\log(\epsilon)/\log(b)$ 越小, 越好. \rightarrow 在资源限制下 (electronic readout, storage, CPU, etc.).

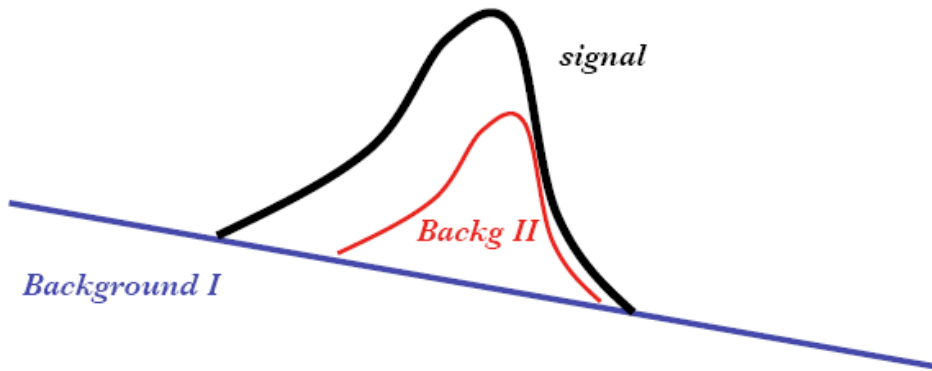


- $\log(\epsilon)/\log(b) < a \rightarrow \epsilon > b^a$ ($0 < \epsilon, b, a \leq 1$).



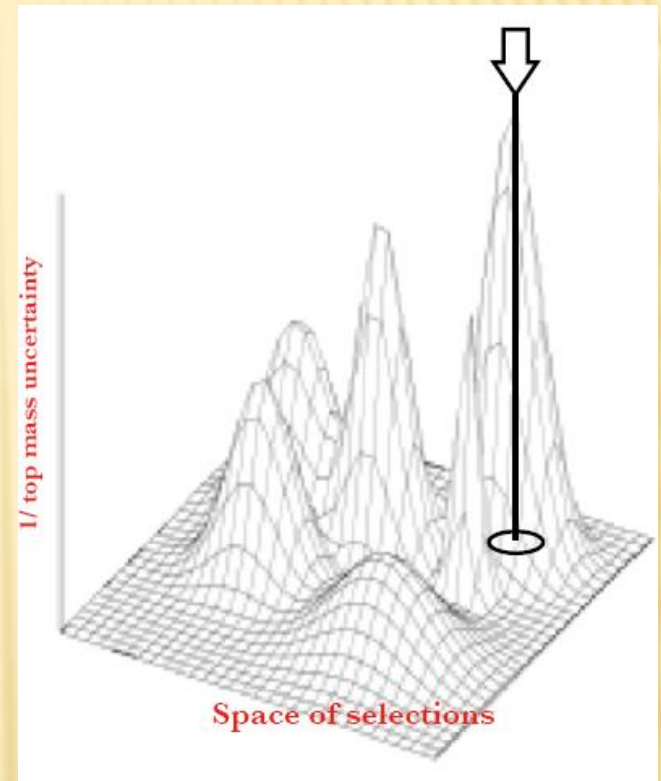
top质量测量中可能人为造成测量误差偏小

A mass measurement depends not just on the **number** of signal or background events but on the **kinematics**



A small number of **background II** can spoil our measurement.

We can tolerate a large number of **background I**

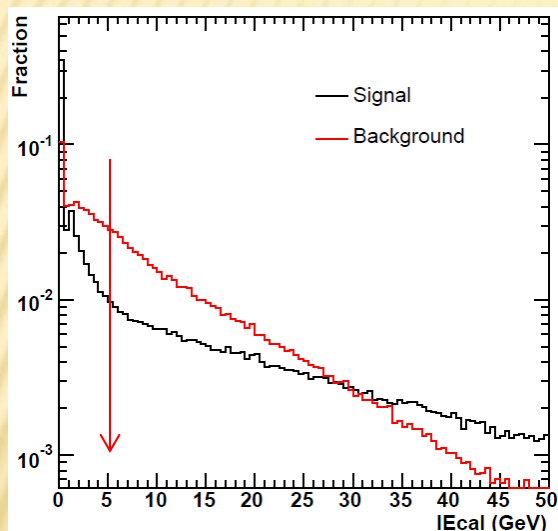


- 将测量误差作为优化标准，当然测量误差人为造成偏小。
- 在对测量结果满意的限制下（比如清晰地看到top信号，对测量误差满意），将 $\log(\epsilon)/\log(b)$ 作为优化标准，就不会有上面的问题。
- 有待检验。

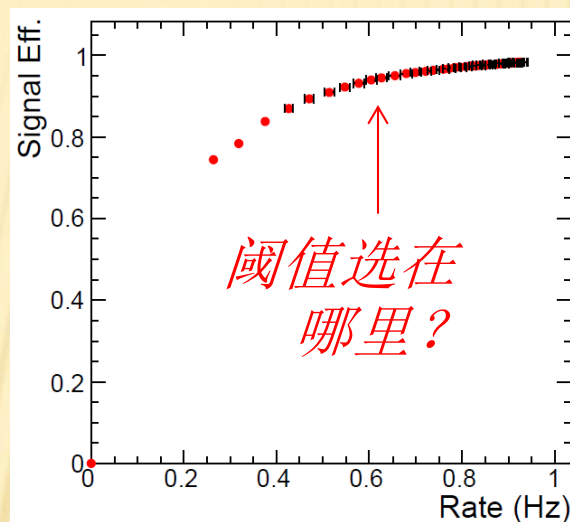
优化的过程

- 对每一个cut, $a_i = \log(\varepsilon_j) / \log(b_j)$, 得到一个最小值 a_i^{\min} ; 假设 $a_1^{\min} < a_2^{\min} < a_3^{\min} < \dots < a_k^{\min}$
 - 如果对效率不满意(太小), 放松 a_k , 去掉 a_k , 放松 a_{k-1} , 去掉 a_{k-1} , ..., 直到我们对效率满意
 - 如果对本底不满意(太多), 将 a_1 从 a_1^{\min} 变化到 a_2^{\min} , 然后将 $a_1 = a_2$ 从 a_2^{\min} 变化到 a_3^{\min} , ..., 直到我们对本底满意.
- → 对于物理分析, 用上面的过程优化选择条件, 直到对MC预测的物理结果满意

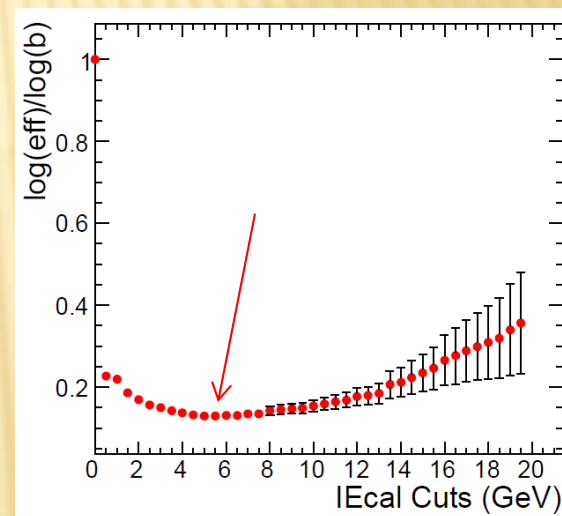
研究一：触发阈值



信号和本底的
触发变量分布



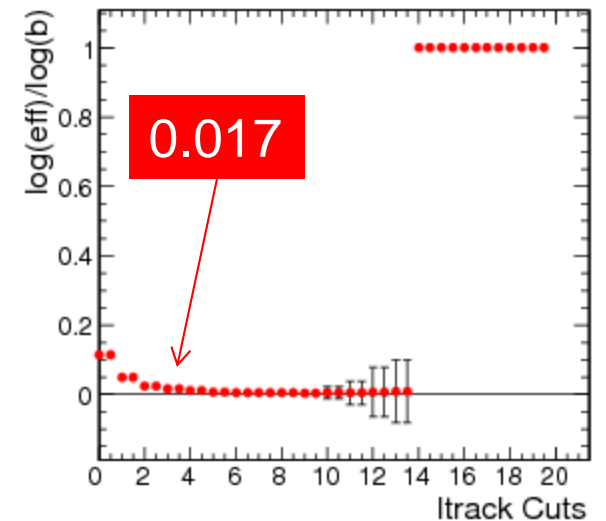
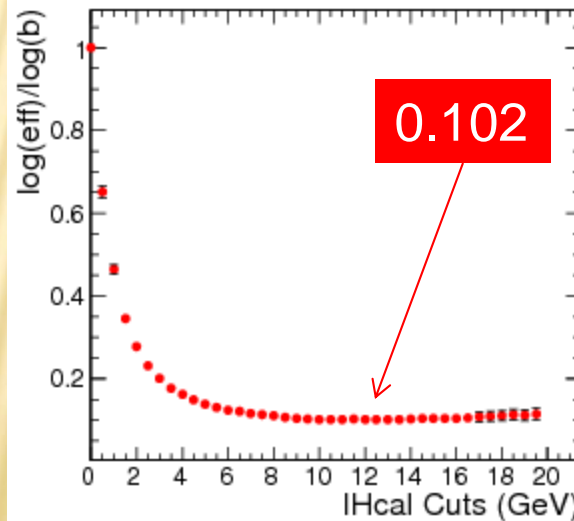
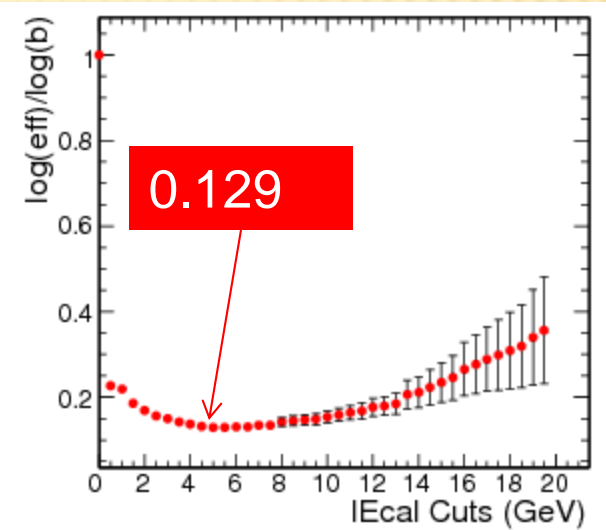
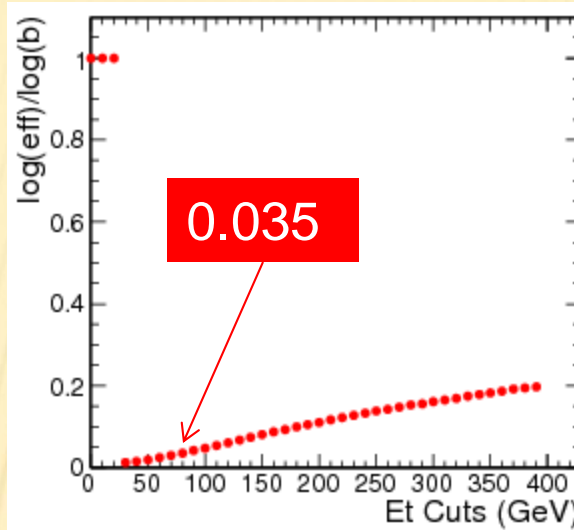
信号效率 vs.
本底效率



$\log(\epsilon)/\log(b)$ vs.
触发阈值

$\log(\epsilon)/\log(b)$ vs. Cuts (原始EM-High-Et)

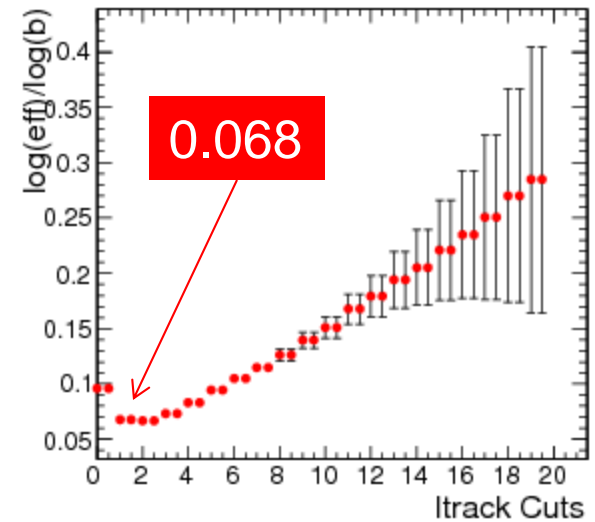
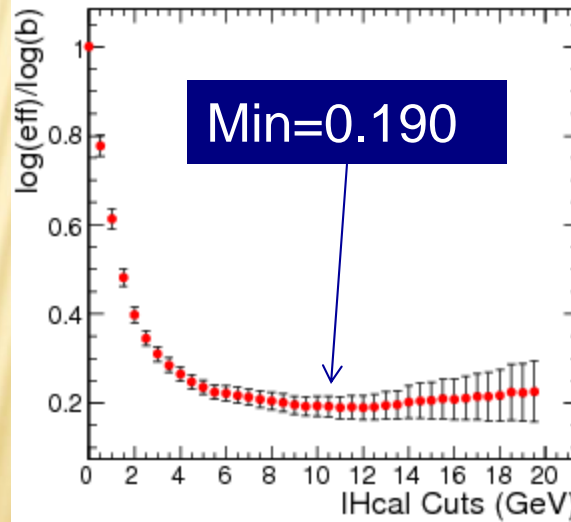
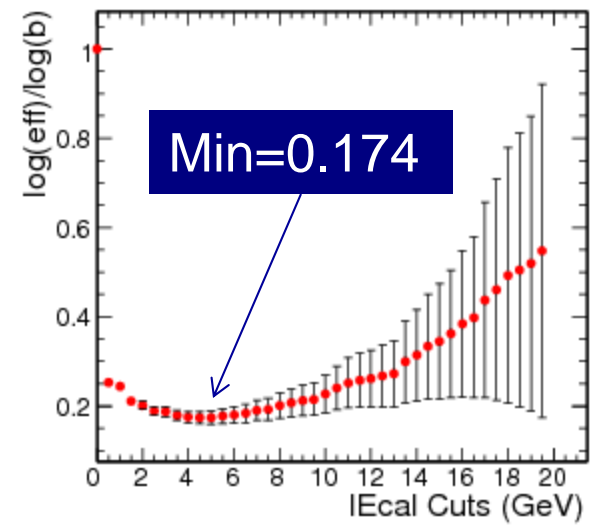
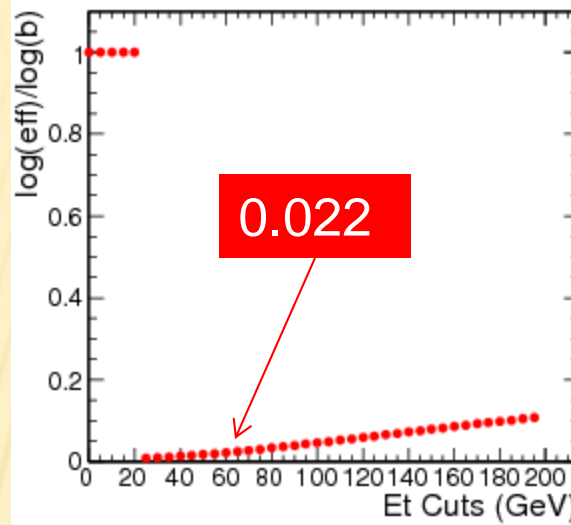
- $E_t > 80$
 $I_{ecal} < 5$
 $I_{hcal} < 12$
 $I_{track} < 4$
- I_{Track} 比
 I_{Ecal} 和 I_{Hcal}
好

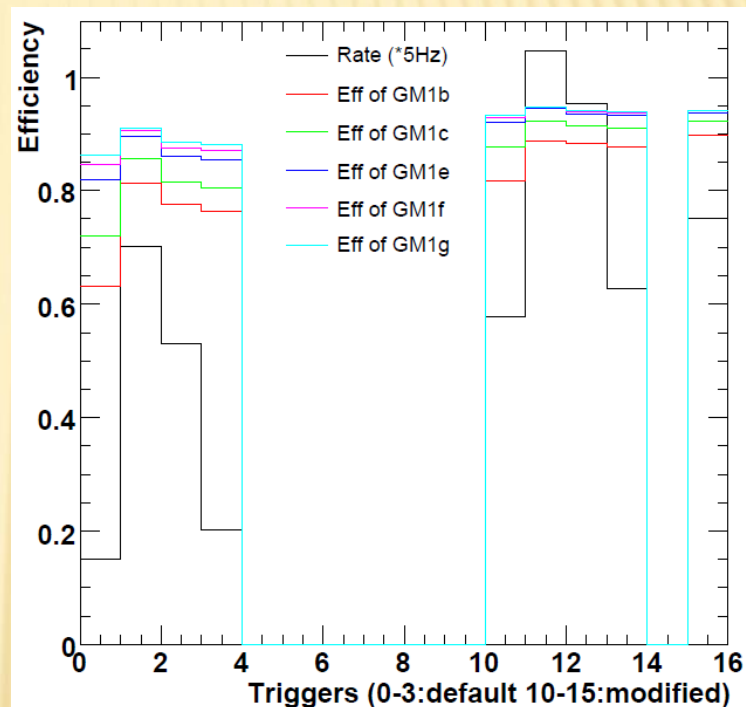
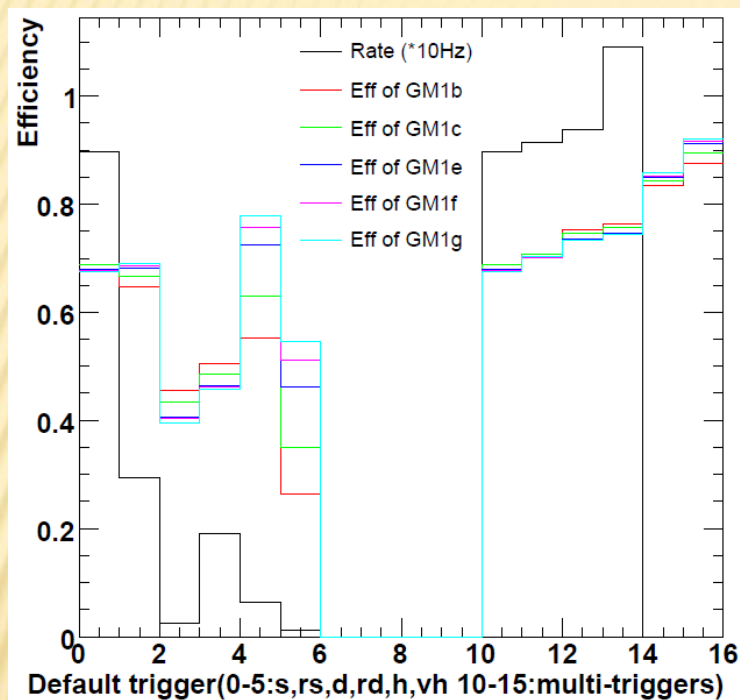


$\log(\epsilon)/\log(b)$ vs. Cuts (优化后EM-High-Et)

- $E_t > 60$;
 $I_{Track} < 2$

- 只有当
 $I_{Track} < 5$
时它才比较
较好

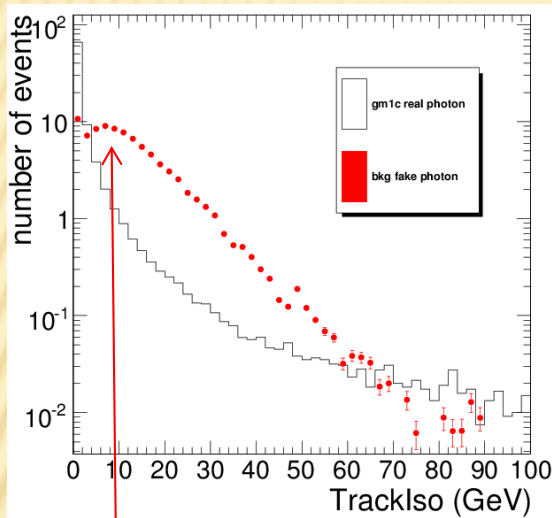




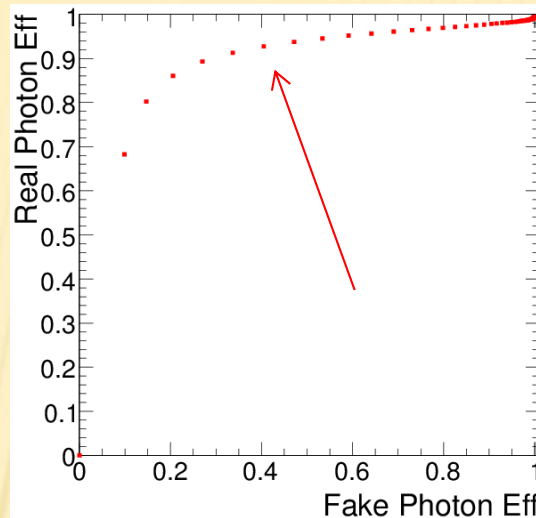
优化前各个触发path的信号效率（彩色直方图）和本底rate率（黑色直方图）

优化后

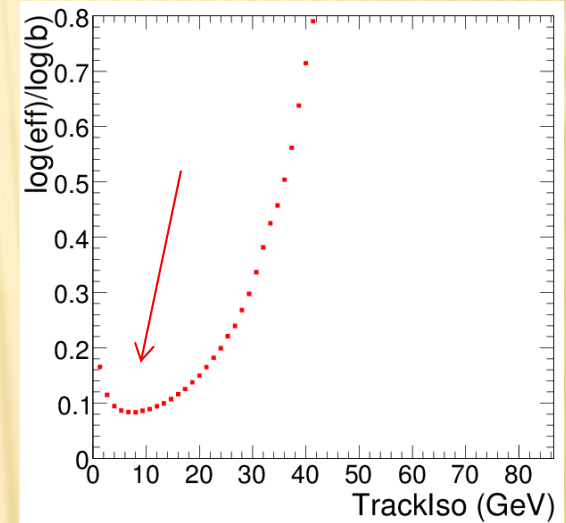
研究二：光子选择



信号和本底的
Cut变量分布



信号效率 vs.
本底效率

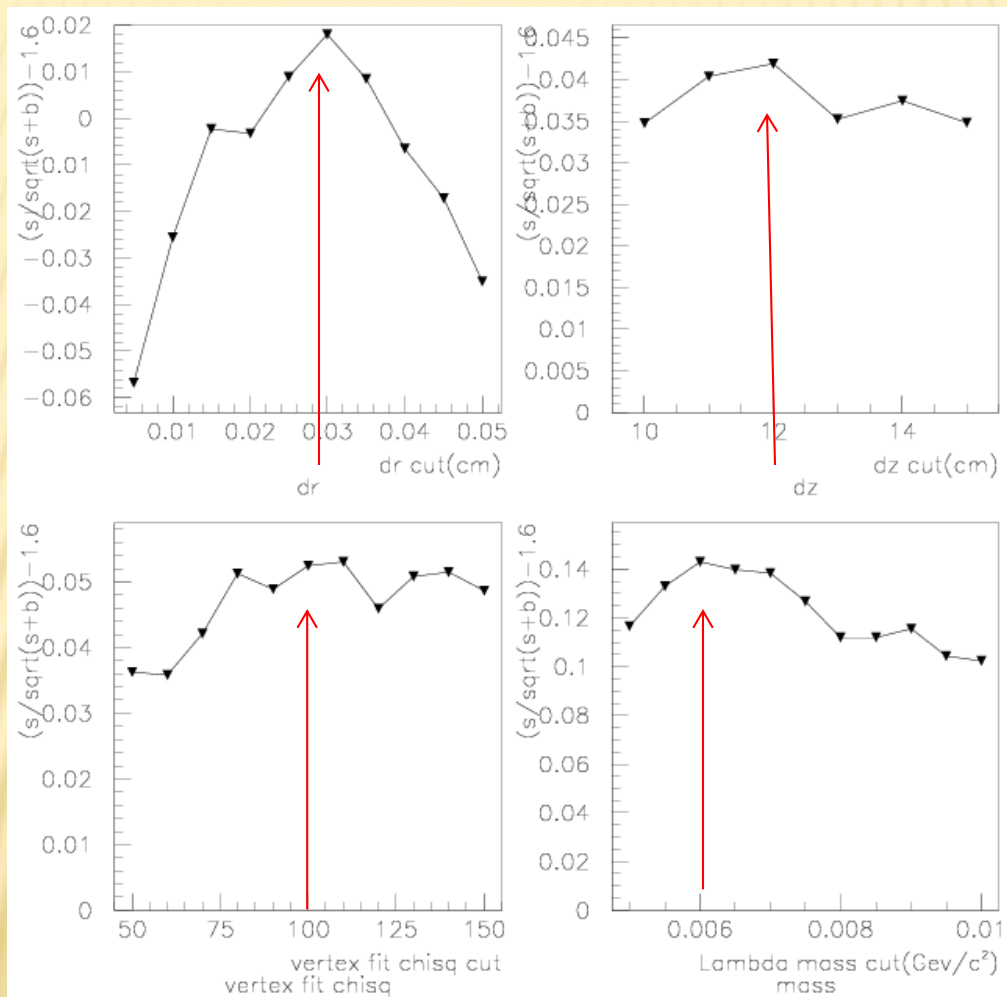


$\log(\epsilon)/\log(b)$ vs.
Cut条件

➤ 优化出的选择条件，被物理分析采用：

- 1) $\text{sigmaEtaEta} < 0.01$
- 2) $\text{HadOverEM} < 0.05$
- 3) $\text{Ecallso}-0.15\text{Pt} < 3\text{GeV}$
- 4) $\text{TrackIso} < 8\text{GeV}$ (上面的图)

$N_S/\sqrt{N_S + N_B}$ 优化选择条件 (tight cuts)



结论

- 提出了一个新的方法，使用信号和本底信息量减少的比值 $\log(\epsilon)/\log(b)$ 做为优化标准，来优化触发、过滤、粒子鉴别、物理分析、多变量分析等等中的阈值设定。
- 方法有信息学和数理逻辑的支撑。
- 方法是客观的，盲分析，减少了人为倾向性。
- 在触发阈值设定和事例选择的研究中被证明是可行的或有用的。
- 在粒子鉴别、物理分析、多变量分析，以及其他领域（比如股市中的收益vs风险），值得尝试应用，预计有所帮助。

谢谢!

Backup Slides

数据

□ 2007 夏季样本, Pythia6, 在14TeV模拟产生

➤ *GMSB* 信号:

- *Five GMSB samples with different Λ parameter*

➤ 本底:

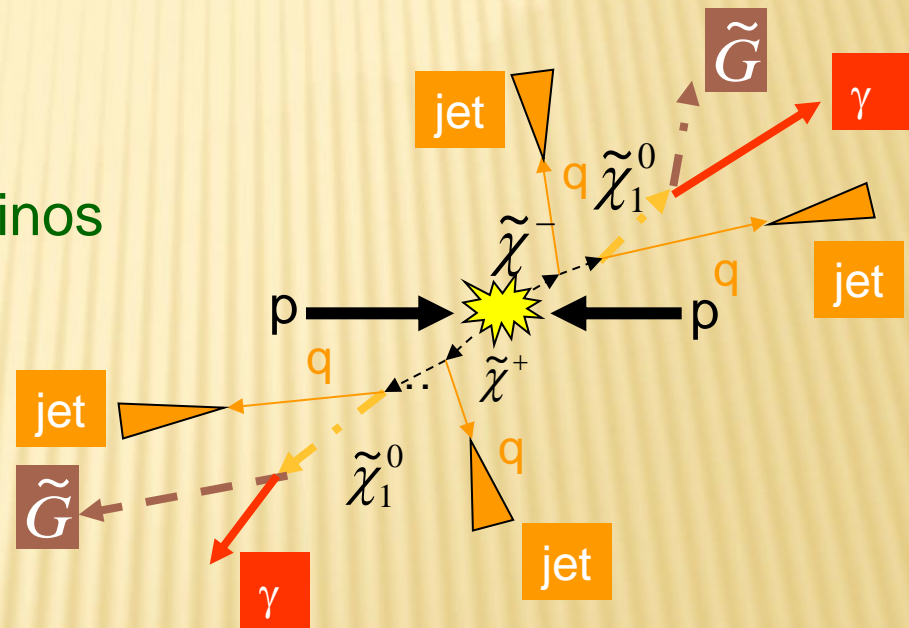
- *Photon+jets (all $pT_{\hat{}}$ bins)*
- *QCD jets (all $pT_{\hat{}}$ bins)*
- *$W \rightarrow e\nu, Z \rightarrow ee$*

□ *CMSSW_1_6_7*; 优化到亮度 $10^{32} \text{ cm}^{-2}\text{s}^{-1}$.

GMSB 信号中的末态光子

- NLSP \rightarrow LSP + photon
- Prompt decay
- Experimental signature
 - high p_T photons
 - large MET due to gravitinos
 - multi-jets

$$\tilde{\chi}_1^0 \rightarrow \tilde{G}\gamma$$



Trigger variables for photons

- *L1Match*: Reconstructed super-cluster in the ECAL is required to match L1 energy deposit in some eta and phi windows.
- *Et*: Et of super-cluster in the ECAL is required to exceed a threshold.
- *IEcal*: ECAL isolation, total Et of all clusters with $\Delta R < 0.3$ around the photon candidate, excluding those belonging to the super-cluster itself.
- *IHcal*: HCAL isolation, total Et of hadron calorimeter towers with $\Delta R < 0.3$ around the photon candidate.
- *ITrack*: Track isolation, number of tracks with $P_t > 1.5 \text{ GeV}$ inside a cone $\Delta R < 0.3$ of photon candidate.

Number of background events processed for rate estimation

QCD Jets	<i>N</i>	QCD Jets	<i>N</i>	Photon Jets	<i>N</i>	Bkg	<i>N</i>
0_15	<i>295,613</i>	_470	<i>88,086</i>	0_15	<i>500,000</i>	Wenu	<i>205,707</i>
15_20	<i>1,255,976</i>	_600	<i>55,000</i>	15_20	<i>509,825</i>	Zee	<i>162,219</i>
20_30	<i>2,513,934</i>	_800	<i>21,974</i>	20_30	<i>606,680</i>		
30_50	<i>2,416,441</i>	_1000	<i>33,330</i>	30_50	<i>510,094</i>		
50_80	<i>2,451,439</i>	_1400	<i>5,299</i>	50_80	<i>169,741</i>		
_120	<i>1,161,823</i>	_1800		_120	<i>164,000</i>		
_170	<i>499,389</i>	_2200		_170	<i>69,993</i>		
_230	<i>428,888</i>	_2600		_300	<i>24,993</i>		
_300	<i>172,619</i>			_500	<i>15,554</i>		
_380	<i>82,998</i>			_7000	<i>6,666</i>		